

# **Statistical Methods to Incorporate External Summary-level Information into a Current Study**

by

Tian Gu

A dissertation submitted in partial fulfillment  
of the requirements for the degree of  
Doctor of Philosophy  
(Biostatistics)  
in the University of Michigan  
2021

Doctoral Committee:

Professor Bhramar Mukherjee, Co-Chair

Professor Jeremy M.G. Taylor, Co-Chair

Assistant Professor Yang Chen

Assistant Professor Peisong Han

Tian Gu

gutian@umich.edu

ORCID iD: 0000-0002-0856-8843

© Tian Gu 2021

## ACKNOWLEDGMENTS

It is unbelievable that I am coming to the end of my graduate school journey. I have had an enjoyable and rewarding six years in Ann Arbor, and I would like to take this opportunity to express my sincere gratitude to many people.

Thanks go first to my Ph.D. co-advisors, Drs. Jeremy M.G. Taylor and Bhramar Mukherjee. I have enjoyed and looked forward to every meeting we had in the past 4 years, in-person or virtually, which happened nearly every week since 2017, rain or shine, even when Dr. Taylor was in New Zealand for a sabbatical or I was in China over the winter breaks. Before the pandemic, I would drop by at Dr. Taylor's office for an extra "office hour" on Friday afternoon, while Bhramar would sometimes invite me to the tea shop for a Sunday afternoon chat, just like friends. In addition to the undoubted A<sup>+</sup> mentoring that I experienced from both of them, I also learned from their unique and charming personalities. Dr. Taylor has been very responsive to all my emails and responsible for every single sentence I wrote in the manuscript. He showed me the most thorough and meticulous attitude, rigorous reasoning, and clear writing in revision after revision until the work met his satisfaction. Bhramar has been brimful of energy and enthusiasm for work, life, and art. Seeing her take on new challenges, balance different roles, accommodate multi-relationships in work, and yet still produce solid work and keep a high spirit in life has influenced me positively. Their knowledgeable minds, humble and curious attitude to research, and passion for novelty provided me support and guidance of dissertation mentorship and how to be a responsible researcher. They will always have my deepest admiration and respect.

My sincere appreciation also goes to Dr. Susan Murray, my advisor for my graduate student research assistant position. Susan has been patiently teaching me every detail in work, turning all tasks into learning opportunities, and recognizing every credit I earned. I truly benefited from her critical thinking in applied statistical analysis and excellent communication skills with non-statisticians. Susan was never sparing with praise and encouragement, and she showed exceptional support and respect towards all my choices. It has been a delightful and rewarding experience working with her and the wonderful collaborators at the pulmonary group in the University of Michigan Medical School: Meilan Han, Wassim Labaki, Margaret Salisbury, Beth Belloli, and Vebha Lama.

Many thanks to members of my committee, Drs. Peisong Han and Yang Chen. Peisong has inspired me to solve the largest obstacle in my second project, a derivation that troubled me for over a half year. He kindly endured my many unannounced visits. Yang has been supportive and responsive to all my requests. I am honored to have had them as my committee members.

In addition to my co-advisors and committee members, I would like to thank professors from the data integration working group—Drs. Vera Baladandayuthapani, Peter Song, Phil Boonstra, Kevin He, Jian Kang, and Hui Jiang, who listened to many of my immature presentations about my preliminary research results and provided constructive insights; especially Dr. Hui Jiang, who guided me on my very first research project in graduate school over the summer of 2016.

I want to thank my lab-mates in both Mukherjee and TaBaBoo lab, especially Dr. Lauren Beesley, who is a master at balancing work and life, and a great friend and mentor who always inspires me academically and mentally; Max Salvatore and Dr. Lars Fritsche, who were the best teammate and mentor I could ask for in the COVID projects; Dr. Phil Boonstra, my academic brother who was kind enough to spent time looking at my R code in my first project and teach me many great R tricks; and Emily Roberts, Elizabeth Chase, Jiacong Du, and Pedro Orozco Del Pino, who have excellent personalities and have been great friends of mine.

I am lucky to have many classmates and friends who supported each other along the way. My graduate school experience would not have been so colorful without the following lovely people: Youfei Yu, Nina Zhou, Ming Tang, Yingchao Zhong, Yiwang Zhou, Zhangcheng Zhao, Jingyi Zhai, Yatong Li, Lan Luo, Jiaqiang Zhu, Yuqing Lin, Boran Gao, and Yancheng Zhao. In addition, I want to thank our school of public health family, the faculty and staff that have made Michigan Biostats feel like home, especially Nicole Fenech, Irene Felicetti, Dr. Kirsten Herold and Dr. Xu Shi.

Shout out to my cutest current office-mates who form the best office in SPH—Lulu Shang, Fatema Shafie, Kim Hochstedler, Tahmeed Tureen, and Jung Yeon Won. Some of my old office-mates and senior peers showed me how to be a Ph.D. student and explore career paths in my early years, including Drs. Jinchunzi Shi, Yumeng Li, Sheng Qiu, Summer Xia, Boxian Wei, Marco Benedetti, Chris Lee and Vincent Tan. Outside school, the friends I made in the Chinese rowing team—Jiayi Tian, Lu Qu, Ting Lin, Tiancheng Zuo, and Yuhan Liang have witnessed a precious time with me between 2018 and 2020.

I want to show my special gratitude to one of my oldest friends, Lv Ni, who has been on my side since middle school and flew from Wisconsin to participate in my master's graduation ceremony. She has always been generous with her time to discuss any questions that I had in my research, even though her schedule is packed with teaching, research, and tedious administrative work as a young assistant professor. She is my go-to person whenever I am stuck in a math problem. I feel safe to have her backing, and she absolutely deserves an independent paragraph in these

acknowledgments.

Finally, I want to thank my parents, Ruijun Wu and Renxu Gu, and my boyfriend, Zizheng Zhang. My parents have always been my role models, giving me so much trust, freedom, and unconditional support. I admire my mom's determination and hard work as a researcher and her courage and vision as a female leader. As a strict father, my dad helped me develop good habits in living and studying when I was little. Although he participated less in my academic life when my college major was completely out of his specialty, he has always been the cheerleader for my mom and me, and he always takes good care of us. Finally, Zizheng has contributed a lot in helping me manage my stress. Although Zizheng has been the biggest distraction of my studying, without his company, I could not have survived the 451-day-long COVID-quarantine working from home between March 14, 2020, and June 8, 2021.

# TABLE OF CONTENTS

ACKNOWLEDGMENTS . . . . .	ii
LIST OF FIGURES . . . . .	vii
LIST OF TABLES . . . . .	ix
LIST OF APPENDICES . . . . .	xi
LIST OF ACRONYMS . . . . .	xii
ABSTRACT . . . . .	xiv

## CHAPTER

<b>1 Introduction . . . . .</b>	<b>1</b>
<b>2 Synthetic Data Method to Incorporate External Information into a Current Study . .</b>	<b>5</b>
2.1 Introduction . . . . .	5
2.2 Method . . . . .	6
2.2.1 General Description of the Problem . . . . .	6
2.2.2 Synthetic Data Method . . . . .	6
2.3 Simulation Study . . . . .	8
2.4 Prostate Cancer Prevention Trial Data Example . . . . .	13
2.5 Algebraic Justification in Two Special Cases . . . . .	15
2.5.1 Estimation and Variance of $\gamma$ . . . . .	15
2.5.2 Description of Two Special Cases . . . . .	17
2.5.3 Summary . . . . .	19
2.5.4 Justification from Another Perspective . . . . .	20
2.6 Discussion . . . . .	22
2.7 Publication . . . . .	23
<b>3 A Meta-Inference Framework to Integrate Multiple External Models into a Current Study . . . . .</b>	<b>24</b>
3.1 Introduction . . . . .	24
3.2 Models and Methods . . . . .	25
3.2.1 General Description of the Problem . . . . .	25
3.2.2 Two Existing Estimators . . . . .	26

3.2.3	Proposed Meta-Framework for Inference . . . . .	27
3.2.4	Asymptotic Normality and Large Sample Results . . . . .	29
3.3	Simulation Studies . . . . .	31
3.3.1	Simulation Settings . . . . .	31
3.3.2	Simulation Results . . . . .	33
3.4	Application to Prostate Cancer Data . . . . .	38
3.5	Discussion . . . . .	40
3.6	Software and Publication . . . . .	44
<b>4</b>	<b>Regression Inference for Multiple Populations by Integrating Summary-Level Data using Stacked Imputations . . . . .</b>	<b>45</b>
4.1	Introduction . . . . .	45
4.2	Models and Methods . . . . .	46
4.2.1	Notation . . . . .	46
4.2.2	Proposed Data Integration and Analysis Strategy . . . . .	48
4.3	Simulation Studies . . . . .	51
4.3.1	Simulation Settings . . . . .	52
4.3.2	Simulation Results . . . . .	54
4.4	Application to Prostate Cancer Data . . . . .	58
4.5	Discussion . . . . .	60
4.6	Software and Publication . . . . .	62
<b>5</b>	<b>Discussion . . . . .</b>	<b>63</b>
	 APPENDICES . . . . .	 68
	 BIBLIOGRAPHY . . . . .	 94

## LIST OF FIGURES

### FIGURE

2.1	Two steps to create the synthetic data . . . . .	7
3.1	Simulation settings snapshot . . . . .	33
3.2	Visualization of the metrics to evaluate the performance in simulations I-IV. Scatter plot shows the covariate-wise relative MSE improvement compared with the direct regression (x axis represents $\hat{\gamma}_X$ ; $\hat{\gamma}_B$ not shown since no external information incorporated; larger y values represent larger MSE loss) while the line plots represent the relative efficiency/SSE/BS improvement compared with the direct regression fitting on a validation dataset of size 1,000 (longer lines represent larger improvement). . . .	35
3.3	Visualization of the metrics to evaluate the performance in simulations V-VI. Scatter plot shows the covariate-wise relative MSE improvement compared with the direct regression (x axis represents $\hat{\gamma}_X$ ; $\hat{\gamma}_B$ not shown since no external information incorporated; larger y values represent larger MSE loss) while the line plots represent the relative efficiency/SSE/BS improvement compared with the direct regression fitting on a validation dataset of size 1,000 (longer lines represent larger improvement). . . .	38
4.1	Diagram of the proposed data integration and analysis strategy. . . . .	49
4.2	Simulation settings snapshot. . . . .	52
4.3	Visualization of simulation I results over increasing synthetic data size (a) point estimates (b) variance estimators vs. the empirical variance (c) different variance estimators of the proposed strategy. . . . .	56
4.4	Visualization of prediction metrics over increasing synthetic data size for simulation II. Larger AUC, smaller SSE and smaller BS represents better prediction. . . . .	57
B.1	Summary of features and assumptions required in the direct regression, CSPML, EB and the proposed approaches. . . . .	78
B.2	Simulation settings of additional scenarios for simulation III in Chapter 3: the joint distribution of (X, B) is misspecified in the external model 1. . . . .	81
B.3	Simulation settings of additional scenarios for simulation IV in Chapter 3: the outcome model is misspecified in the external model 3. . . . .	82
B.4	Simulation results of additional scenarios of simulation III in Chapter 3. . . . .	83
B.5	Simulation results of additional scenarios of simulation IV in Chapter 3. . . . .	84
C.1	Results of Simulation C.2.1 over increasing synthetic data size (a) point estimates (b)variance estimation vs. Monte Carlo empirical variance (c) different variance estimators of the proposed method. . . . .	89



C.2	Results of Simulation C.2.2 over increasing synthetic data size (a) point estimates (b) different variance estimators of the proposed method. . . . .	90
C.3	Results of Simulation C.2.3 over increasing synthetic data size (a) point estimates (b) different variance estimators of the proposed method. . . . .	91
C.4	Results of Simulation C.2.4.1 over increasing synthetic data size (a) point estimates (b) different variance estimators of the proposed method. . . . .	92
C.5	Results of Simulation C.2.4.2 over increasing synthetic data size (a) point estimates (b) different variance estimators of the proposed method. . . . .	93

## LIST OF TABLES

### TABLE

2.1	Simulation results for scenario 1 with Gaussian Y, one Gaussian B and nine correlated X's. For each method, we report mean (SD) [95 % coverage rate] and MSE across 500 simulated datasets. . . . .	10
2.2	Simulation results for scenario 2 with binary Y, one Gaussian B and nine correlated X's. For each method, we report mean (SD) [95 % coverage rate], average scaled Brier score and AUC across 500 simulated datasets . . . . .	11
2.3	Simulation results for scenario 3 with binary Y, one binary B and nine correlated X's. For each method, we report mean (SD) [95 % coverage rate], average scaled Brier score and AUC across 500 simulated datasets . . . . .	12
2.4	Simulation results for scenario 4 with binary Y, binary B1 and continuous B2 and nine correlated X's. For each method, we report mean (SD) [95 % coverage rate], average scaled Brier score and AUC across 500 simulated datasets . . . . .	13
2.5	Results of the expanded PCPThg model. For each method, point estimate (standard error) from the internal dataset, and the scaled Brier score and the AUC from the validation dataset. The sample size of the internal dataset is 679. The sample size of the validation dataset is 1,218. Replicate number S=10 giving m=6,790 in the synthetic data method. . . . .	14
2.6	Summary of the special case 1 (Y and B are Gaussian) . . . . .	19
2.7	Summary of the special case 2 (Y, X and B are binary) . . . . .	19
3.1	Results of simulation I–IV. Internal dataset had size n=200; green represents good performance with small bias and large efficiency gain, yellow represents underestimated ESE (in square brackets) compared with SD (in round brackets), and red represents poor performance of bias/95% coverage rate. . . . .	34
3.2	Results of simulation V–VI. Internal dataset had size n=500; green represents good performance with small bias and large efficiency gain. . . . .	37
3.3	Results of the real data example for predicting the risk of high-grade prostate cancer. Internal dataset of size n=678; validation dataset of size $N_{\text{test}}=1,174$ ; ↓ %, percentage of ESE decrease compared with the direct regression; Firth correction applied in direct regression, and estimated PCPThg and ERSPC; REF, reference; green represents good performance with small bias and large efficiency gain, and grey represents no efficiency improvement due to no added information from the external calculator. . . . .	41

4.1	Results of the data example for predicting the risk of high-grade prostate cancer. Internal dataset of size $n=678$ ; validation dataset of size $N_{\text{test}}=1,174$ ; SE, standard error; green represents good performance with large precision gain or better model calibration, grey represents certain predictor was not used in the external calculator, yellow represents population-specific effect different from the internal population, and red represents poor performance compared with direct regression; SE in the proposed method are (StackImpute SE) [bootstrap SE from 500 replicates]. . . . .	59
A.1	Formulas of $\gamma$ in terms of $P(B XY)$ and $P(XY)$ . . . . .	73

**LIST OF APPENDICES**

**A Appendix of Chapter 2 . . . . . 68**

**B Appendix of Chapter 3 . . . . . 77**

**C Appendix of Chapter 4 . . . . . 86**

## LIST OF ACRONYMS

**ARE** asymptotic relative efficiency

**AUC** area under the curve

**BS** scaled Brier score

**CML** constrained maximum likelihood

**CSPML** constrained semi-parametric maximum likelihood

**DRE** digital rectal exam

**EB** empirical Bayes

**ERSPC** European Randomized Study of Screening for Prostate Cancer risk calculator 3

**ESE** estimated standard error

**IVW** inverse variance weighted estimator

**MLE** maximum likelihood estimation

**MSE** mean squared error

**OCWE** optimal covariates weighted estimator

**PCPT<sub>hg</sub>** Prostate Cancer Prevention Trial (high-grade prostate cancer risk calculator)

**PCA3** biomarker prostate cancer antigen 3

**PSA** prostate-specific antigen

**SC-Learner** selective coefficient learner

**SD** Monte Carlo standard deviation from 500 simulations

**SSE** sum of squared error

**TRUS-PV** transrectal ultrasound prostate volume

**T2:ERG** biomarker TMPRSS2:ERG

## ABSTRACT

In the era of big data, it is becoming increasingly common for researchers to consider incorporating external information from large studies to improve the accuracy of statistical inference instead of relying on a modestly sized dataset collected internally. We consider a general statistical problem where there are some known regression models or risk calculators to predict an outcome of interest from a set of commonly used predictors. Different types of summary information are available for these external models. An internal modest-sized dataset containing individual-level data for the variables in the known models and some new variables is available for our current analysis. In all three chapters below, we consider different settings to achieve the same goal—to build an improved prediction model that includes the new variables, using both the internal individual-level data and summary information obtained from the known external model(s).

In Chapter 2, we focus on the simple case where there is only one large, well-characterized previous study from the external population. We propose a synthetic data approach, which first converts the external information into synthetic data, and then analyzes a combined dataset consisting of the observed internal data and the synthetic data. A theoretical justification and extensive simulation studies establish the efficiency gain and improved prediction performance of the proposed data integration method. We also illustrate that even under less restrictive requirements on the information that is available externally, the combined estimates have the same asymptotic properties as an alternative constraint maximum likelihood estimation approach.

In Chapter 3, we consider a more complicated but quite plausible situation where several external prediction models are available to aid inference and prediction for the internal study. We assume that each of the external studies developed a prediction model for the same outcome but may use a slightly different set of covariates. We propose a meta-inference framework using an empirical Bayes estimation approach, which adaptively combines the estimates from the external models. This adaptive approach diminishes the influence of information that is less compatible with the internal data while balancing the bias-variance trade-off. The estimators we proposed are more efficient than the naive analysis of the internal data.

In Chapter 4, we first extend the synthetic data method from Chapter 2 to accommodate the situation with multiple external prediction models, and further allow for heterogeneity of covariate effects across the external populations. Each external model could potentially be built on slightly

different subsets of covariates that are measured in the internal study. The proposed approach generates synthetic outcome data in each population, uses stacked multiple imputation to create a long dataset with complete covariate information, and finally analyzes the imputed data with weighted regression. Leveraging multiple sources of auxiliary information from a broad class of externally fitted predictive models or established risk calculators based on parametric regression or machine learning methods, this new strategy can make statistical inference more accurate for both the internal population and the external populations.

We evaluate the proposed methods through extensive simulations and apply them to improve models for predicting the risk of high-grade prostate cancer.



# CHAPTER 1

## Introduction

Increasingly, researchers are considering incorporating external information from large-scale studies to improve statistical inference rather than using the limited-sized data that are available to each investigator. Examples of this are borrowing strength from historical control data to leverage the treatment effect in small-sample clinical trials [Viele et al., 2014, Dejardin et al., 2018, Li and Song, 2020], combining separate probability samples [Bycroft, 2011, Yang and Kim, 2020] and incorporating external data sources for improved causal inference [Yang and Ding, 2020a]. However, challenges exist, such as data sharing, storage, and privacy issues to access publicly available individual-level large data, so often only the summary information is reported. Examples of such data sources include publications, online risk calculators, census data and population-based biobank data. Therefore, general frameworks that integrate the individual-level data and the summary-level external information are needed.

As a motivating example, it is common in predictive modeling that researchers want to include new predictors to update the traditional risk models in clinical biomedicine, such as adding genetic risk variants and mammographic density to the breast cancer risk calculator [Gail et al., 1989]. These traditional models are usually constructed from large datasets using principled statistical methods to predict a measure of risk or disease state, treating the patient characteristics as predictors. The patient characteristics, denoted by  $\mathbf{X}$ , can range from traditional epidemiologic, clinical and behavioral variables to well-known imaging, genetic and other molecular biomarkers. The predicted outcome variable  $Y$ , and the predictors  $\mathbf{X}$ , are often assumed to be connected through a regression model of the form  $Y|\mathbf{X}$ . The individual-level original data that were used to construct this model are usually not available to the public, but what are accessible are certain forms of summary-level information. This information can be available in the form of coefficient estimates for the fitted model [Thompson et al., 2006, Roobol et al., 2012] or individual prediction probabilities, of which the underlying model may or may not be known, especially when the external information comes in the form of a “black box” algorithm, i.e. an algorithm that provides a predicted probability, but the underlying model is not necessarily simple or transparent or even known

[Stephan et al., 2003, Osareh and Shadgar, 2010, Estiri et al., 2021]. Furthermore, different external models predicting the same outcome may use different sets of predictors [Thompson et al., 2006, Roobol et al., 2012].

While these existing models are often based on traditional epidemiologic and behavioral risk factors and well-established biomarkers, wider availability of high throughput data and novel assay technologies are generating new candidate biomarkers, say  $B$ , for possible inclusion in existing risk prediction models. Due to the potential improvement of prediction accuracy of the current model, it is ideal to incorporate  $B$  into the well-established model  $Y|X$ , and construct an expanded prediction model of interest  $Y|X, B$ . However, it is very likely that  $B$  and  $X$  are assessed only on participants in a study of moderate size and cannot be retrospectively measured on the much larger population used for  $Y|X$  model. It is natural to consider using the information from the well-established model to increase the accuracy of the expanded model. This represents a general statistical challenge to build a good model for  $Y|X, B$  that uses both the known external information from the  $Y|X$  model and the individual-level data from a small sample dataset of  $Y$ ,  $X$  and  $B$ .

There exist proposals in the literature to incorporate external information into regression estimation. Imbens and Lancaster [1994] investigated how aggregate data (e.g. the population average of the response) could be used to improve maximum likelihood estimates in a regression model. Grill et al. [2015] proposed a simple method of incorporating new markers into an existing calculator via Bayes Theorem. Studies on this topic begin with incorporating external auxiliary information from large data, such as census or population-based biobank data, to improve the statistical inference of the internal study, assuming the model that relates the variables is fully or partially shared between data sources. The topic of whether distributions are similar between populations is called transportability [Bareinboim and Pearl, 2013] and is crucial to consider in the field of data integration. Qin [2000], Han and Lawless [2019] proposed possible solutions using empirical likelihood, Chatterjee et al. [2016], Cheng et al. [2018] demonstrated it from the perspective of constrained maximum likelihood while Cheng et al. [2019] proposed to use the Bayesian approach. The performance of various approaches was assessed in a simulation study by Grill et al. [2017]. Estes, Mukherjee, and Taylor [2017] then relaxed the transportability assumption of Chatterjee et al. [2016]’s by constructing an empirical Bayes estimator that protected against the potential bias.

Recent studies started to tackle challenges such as accommodating multiple external data sources and/or heterogeneity exists among data sources. Kundu, Tang, and Chatterjee [2019] extended the work of Chatterjee et al. [2016] to a meta-analysis setting. Chen et al. [2020] harmonized the difference of aggregate information among data sources through a penalty function while Yang and Ding [2020b] employed a sensitivity parameter to quantify such systematic differences. Moreover, there could be other sources of information variation across the models. For example, different external studies may use a different subset of covariates and the underlying prediction

model may be parametric or constructed by machine learning approaches. The summary-level information may contain estimated regression coefficients or fitted predictions.

In this dissertation, we consider a general statistical problem where there are some known regression models or risk calculators to predict an outcome of interest from a set of commonly used predictors. Different types of summary information are available for these external models. An internal modest-sized dataset containing individual-level data for the variables in the known models as well as some new variables is available for our current analysis. In all three chapters below, we consider different settings aiming to achieve the same goal—to build an improved prediction model that includes the new variables, using both the internal individual-level data and summary information obtained from the known external model(s).

In Chapter 2, we focus on the simple case where there is only one large, well-characterized previous study from the external population. We propose a synthetic data approach, which first converts the external information into synthetic data, and then analyzes a combined dataset consisting of the observed internal data and the synthetic data. A theoretical justification and extensive simulation studies establish the efficiency gain and improved prediction performance of the proposed data integration method. We also illustrate that even under less restrictive requirements on the information that is available externally, the combined estimates have the same asymptotic properties as an alternative constraint maximum likelihood estimation approach.

In Chapter 3, we consider a more complicated but quite plausible situation where several external prediction models are available to aid inference and prediction for the internal study. We assume that each of the external studies developed a prediction model for the same outcome but may use a slightly different set of covariates. We propose a meta-inference framework using an empirical Bayes estimation approach, which adaptively combines the estimates from the external models. This adaptive approach diminishes the influence of information that is less compatible with the internal data while balancing the bias-variance trade-off. The estimators we proposed are more efficient than the naive analysis of the internal data.

In Chapter 4, we first extend the synthetic data method from Chapter 2 to accommodate the situation with multiple external prediction models, and further allow for heterogeneity of covariate effects across the external populations. Each external model could potentially be built on slightly different subsets of covariates that are measured in the internal study. The proposed approach generates synthetic outcome data in each population, uses stacked multiple imputation to create a long dataset with complete covariate information, and finally analyzes the imputed data with weighted regression. Leveraging multiple sources of auxiliary information from a broad class of externally fitted predictive models or established risk calculators based on parametric regression or machine learning methods, this new strategy can make statistical inference more accurate for both the internal population and the external populations.

In all the chapters, we evaluate the proposed methods through extensive simulations and apply them to improve models for predicting the risk of high-grade prostate cancer.

## CHAPTER 2

# Synthetic Data Method to Incorporate External Information into a Current Study

### 2.1 Introduction

In Chapter 1, we introduced some of the existing methods to solve this genre of problem. In general, the constrained maximum likelihood (CML) approaches require a specific form for the external information, e.g. estimated coefficients from a correctly specified mean model and assumptions regarding the transportability of the distribution of  $Y$ ,  $X$ ,  $B$  across the internal and external sample. The CML approach proposed in Cheng et al. [2018] also requires the specification of a model for  $B|X$  and relies on some parametric assumptions. Chatterjee et al. [2016] described a constrained semi-parametric maximum likelihood (CSPML) method by converting the external summary-level information into a constraint and then maximizing the internal data likelihood subject to this constraint. Although this CSPML approach does not require the  $Y|X$  model to be correctly specified or a model for  $B|X$ , it requires the transportability of the joint distribution of  $Y$ ,  $X$ ,  $B$ . Estes, Mukherjee, and Taylor [2017] and Cheng et al. [2018] have found that violation of this assumption and the small sample size in the internal data will cause sensitive and unstable estimation.

In this chapter, we propose a synthetic data framework as a more flexible alternative solution to the CML approach, motivated by methods developed in the survey methodology literature [Raghu-nathan et al., 2003, Reiter and Kinney, 2012, Reiter, 2002]. In this approach, synthetic data for  $X$  and  $Y$  are generated from the observed data and the  $Y|X$  model, respectively, and then added to the observed data, then from this combined dataset a model for  $Y|X, B$  is built. Our method relaxes the requirement on the information that is available from the external model such that the only requirement is the ability to generate predictions of  $Y$  given  $X$ , which can make use of external information that comes in the form of a “black box” algorithm, i.e., an algorithm that provides a predicted probability, but the underlying model is not necessarily simple or transparent or even known.

The following is the structure of the remainder of this chapter: in Section 2.2, we introduce the notation, assumptions and implementation of the proposed synthetic data method. In Section 2.3, under various simulation scenarios, we evaluate the performance of the synthetic data method. We demonstrate the proposed method through an application to the Prostate Cancer Prevention Trial data in Section 2.4. We provide some theoretical justification and insight for the synthetic data method in Section 2.5. In two special cases we show that with a very large number of synthetic observations, our approach gives identical asymptotic variances for the parameters of the  $Y|X, B$  model as the CML estimation approach that exists in the literature. Because the CML is a maximum likelihood estimator, it is optimal if the models are correctly specified. Since the synthetic data method has the same asymptotic variance, it can also be considered optimal. Concluding remarks are presented in Section 2.6.

## 2.2 Method

### 2.2.1 General Description of the Problem

Let  $Y$  denote the outcome of interest, which can be either continuous or binary.  $X$  is a set of  $p$  standard variables and let  $B$  denote a new biomarker. There are two populations, an external population for which we do not have individual-level data and an internal population for which we do have a dataset of size  $n$  with subject-level data. We will assume that the distributions of  $Y|X, B$  are the same in the two populations, and likewise for the  $Y|X$  distribution. Our target of interest is the mean structure of  $Y|X, B$ :

$$g[E(Y|X, B)] = \gamma_0 + \gamma_{X_1}X_1 + \dots + \gamma_{X_p}X_p + \gamma_B B, \quad (2.1)$$

where  $g$  is the known link function. We assume that a small dataset of size  $n$  with variables  $Y, X$  and a new covariate  $B$  is available to us for building the model of interest.

We assume a large, well-characterized previous study from the external population describes the provided information on the calculated distribution of  $Y|X$ . This information can come in various forms, including partial or full knowledge of the distribution of the  $Y|X$  model.

### 2.2.2 Synthetic Data Method

We propose an algorithmic approach that can produce synthetic data on  $(Y, X, B)$ , by using the combination of the available information from the established model and the observations from the current dataset. The synthetic data would incorporate the external information as well as enlarge

the sample size, and thus it helps improve the inference about coefficients  $\gamma$  in model 2.1, compared to just analyzing the small dataset based on the observed data.

The synthetic data approach consists of creating  $m$  additional synthetic data observations, and then analyzing the combined dataset of size  $n+m$  to estimate the parameters of model 2.1. The synthetic data are created in two steps as shown in Figure 2.1. In step 1, we replicate  $\mathbf{X}$  a large number (say  $S$ ) times in blocks of  $n$  rows to create  $m = nS$  additional records. In step 2, we generate pseudo data called  $Y^*$  from the known  $Y|\mathbf{X}$  distribution for these new  $m$  records. Finally, we combine the synthetic observations with the original dataset, and we note that the combined data will now have missing values of  $B$  for  $m$  observations. The combined data is then analyzed to give an estimate of  $\gamma$ .

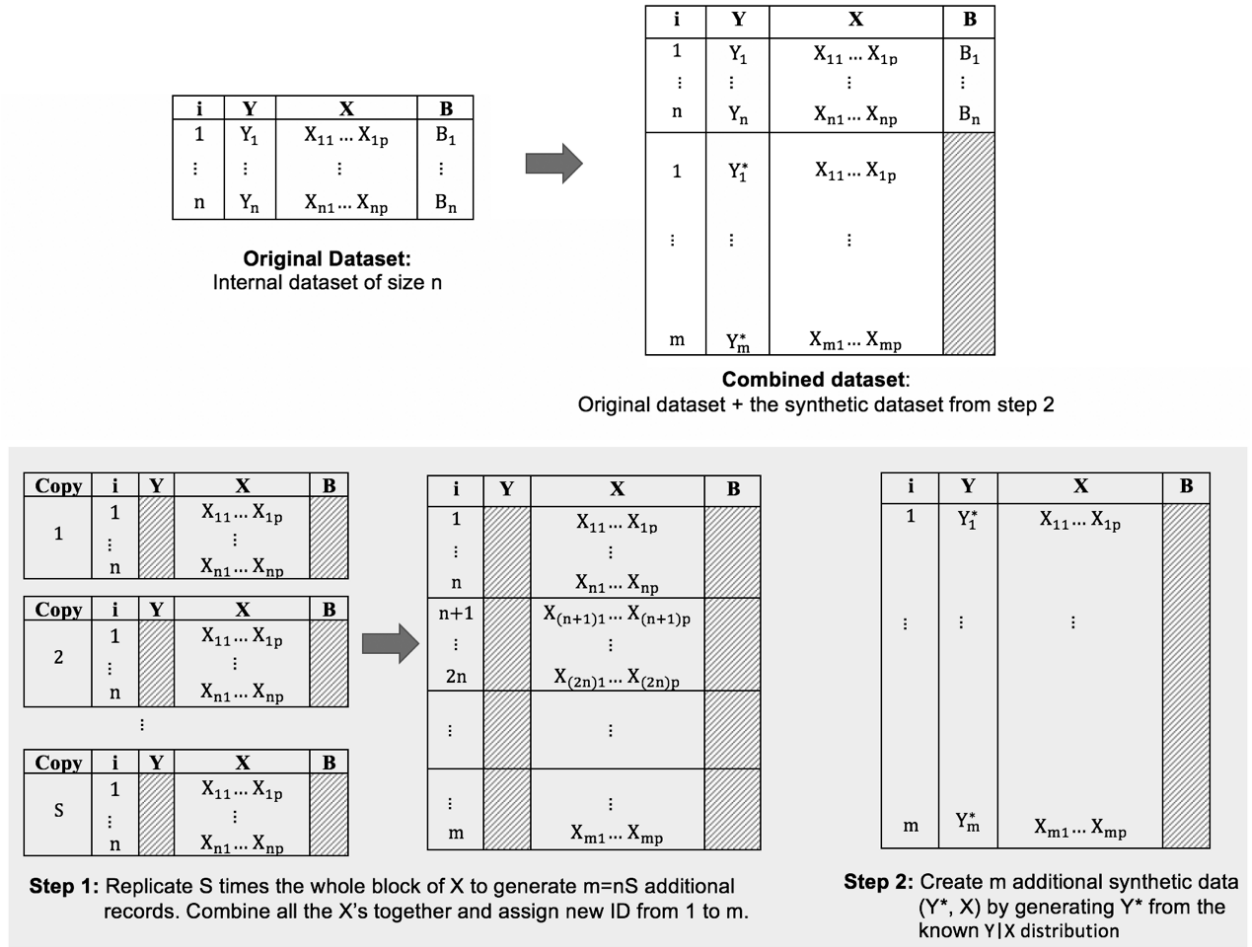


Figure 2.1: Two steps to create the synthetic data

There may be different ways in which the combined data can be analyzed. In Section 5 we present two special cases for which a closed-form MLE of  $\gamma$  exists for the combined dataset of size  $n+m$ . In other cases, like the simulation study settings in Section 3, no closed-form solution

for the MLE of  $\gamma$  exists, and our proposed approach to deal with missing data is to use multiple imputation to impute the  $m$  missing values of  $B$ . Multiple imputation is a general procedure for analyzing datasets with missing values. It consists of defining a procedure to fill in the missing values, then applying that procedure many times to create many separate complete datasets. Each completed dataset is then analyzed and the results of these separate analyses are combined to give final estimates. In this particular case the multiple imputation approach requires us to specify a parametric model  $(B|X, Y)$ , from which we draw 50 values of  $B$  to give 50 completed datasets. Then we fit model 2.1 for each complete data  $(Y, X, B)$  of size  $n+m$ . We then average the estimates of  $\gamma$  from the 50 complete datasets, and compute the total variance using Rubin’s rules [Rubin, 1987]. We then proceed with inference.

Multiple imputation has the additional advantage of being able to handle multiple biomarkers in  $B$ , some of which may be discrete and some continuous. It also allows for flexible structure for the conditional mean model for each biomarker in  $B$  given all other variables in the dataset, such as the possibility to incorporate non-linearity and interactions. For implementing multiple imputation, we use the R package MICE [Van Buuren and Oudshoorn, 2011]. We use the function *mice* with imputation algorithm *logreg* (the Bayesian logistic regression model with flat prior) for the imputation of a binary  $B$  and the imputation algorithm *norm* (the Bayesian linear regression model) for the imputation of a continuous  $B$ . In the situation in which there are multiple  $B$ ’s, say  $B_1$  and  $B_2$ , imputations are done sequentially. That is, first draw  $B_1$  from the  $B_1|X, Y, B_2$  distribution, then draw  $B_2$  from the  $B_2|X, Y, B_1$  distribution, and iterate between  $B_1$  and  $B_2$ .

## 2.3 Simulation Study

To assess the performance of the proposed synthetic data method for both estimation and prediction, we conduct simulation studies under four different scenarios. Each scenario has a different true distribution for  $Y|X, B$  and for  $B|X$  for the internal data. For both the outcome and the predictors, we consider both continuous and binary variables to illustrate the computational implementation in a range of situations. We also consider the situation of multiple  $B$ ’s to evaluate the applicability of the synthetic method in multi-dimensional cases. In some cases a misspecified imputation model is used within the synthetic data approach, thus allowing us to evaluate the robustness of the method. Only in special cases (see Section 2.5) can we provide a theoretical justification for the synthetic data approach, thus the simulations are intended to provide numerical properties of the synthetic data approach in situations where the relevant theoretical properties are not yet available.

In real situations, we expect a moderate number of  $X$  variables, and their joint distribution could be quite complex with skew distributions and correlations between different  $X$ ’s. To achieve this



we adopted a procedure of generating  $\mathbf{X}$ 's as described in Xu, Daniels, and Winterstein [2016]. We generate 9 correlated  $\mathbf{X}$ 's in each of the four simulation scenarios as described below:

$$u_j \sim N(0, 1), \quad j=1 \dots 5, \quad X_1 \sim N(0, 1)$$

$$X_j | u_1, \dots, u_5 = \begin{cases} u_1 u_j + \epsilon_j, \epsilon_j \sim \frac{2}{3}N_+(0, 0.2) + \frac{1}{3}N_-(0, 0.1) & j = 2, \dots, 5 \\ u_2 u_{j-3} + \epsilon_j, \epsilon_j \sim \frac{1}{4}N_+(0, 0.1) + \frac{3}{4}N_-(0, 0.3) & j = 6, \dots, 8 \\ u_3 u_{j-5} + \epsilon_j, \epsilon_j \sim \frac{4}{5}N_+(0, 0.4) + \frac{1}{5}N_-(0, 0.1) & j = 9 \end{cases}$$

where  $N_+$  and  $N_-$  represent half normal distributions, and either one or the other is selected with the shown probability. We then generate  $\mathbf{B}$  from the  $\mathbf{B}|\mathbf{X}$  distribution, and finally generate  $\mathbf{Y}$  from the  $\mathbf{Y}|\mathbf{X}, \mathbf{B}$  distribution. For scenario 1 (where  $\mathbf{Y}$  is continuous, as described below) the form of the external model for  $\mathbf{Y}|\mathbf{X}$  is readily available. For cases (scenarios 2, 3, and 4 where  $\mathbf{Y}$  is binary, as described below) where the closed-form of model  $\mathbf{Y}|\mathbf{X}$  is not available, we numerically derive the external model  $\mathbf{Y}|\mathbf{X}$ . Specifically, we generate an independent dataset of  $(\mathbf{Y}, \mathbf{X}, \mathbf{B})$  of size 10000 and fit a linear or logistic regression model  $g[E(\mathbf{Y}|\mathbf{X})]$  depending on the type of  $\mathbf{Y}$ . The estimated coefficients of this model serve as the external information we obtained from the established model  $\mathbf{Y}|\mathbf{X}$ .

For each simulation scenario, we first simulate 500 datasets of size  $n$ . Then we create the synthetic data following the steps introduced in Section 2, and combine them with the original data to get 500 datasets of size  $n+m$  with  $m$  missing  $\mathbf{B}$  values. For each simulated dataset, we create 50 complete datasets by imputing the missing  $\mathbf{B}$  values given  $\mathbf{Y}$  and  $\mathbf{X}$ . In all four scenarios, we use linear additive models for imputing from the  $\mathbf{B}|\mathbf{X}, \mathbf{Y}$  distribution, without including any interaction terms. We compare the results of the synthetic data method to the direct MLE, which uses the complete dataset of size  $n$ , in terms of estimation accuracy and prediction ability. We report the average estimated coefficients, standard deviation and 95% coverage rate for  $\hat{\gamma}$ . To measure the predictive performance, we generate a new dataset of size 1500 for each scenario, and evaluate the prediction  $\hat{Y}_i$  on this new dataset. In the new validation dataset, let  $\hat{p}$  or  $\bar{Y}$  denote the average of the generated  $\mathbf{Y}$  values. For the continuous  $\mathbf{Y}$ , we use the mean squared error (MSE) defined as  $\sum_{i=1}^{1500} (\hat{Y}_i - \bar{Y})^2 / \sum_{i=1}^{1500} (Y_i - \bar{Y})^2$ . For binary  $\mathbf{Y}$ , we use AUC and scaled Brier score, defined as  $\sum_{i=1}^{1500} (Y_i - \hat{Y}_i)^2 / \sum_{i=1}^{1500} (Y_i - \hat{p})^2$ , as measures of discrimination and calibration.

The four simulation scenarios and results are described as follows:

- **Scenario 1:**  $\mathbf{Y}$  and  $\mathbf{B}$  are Gaussian distributed. The true model of  $\mathbf{Y}|\mathbf{X}, \mathbf{B}$  is  $Y_i | \mathbf{X}_i, B_i = 0.5 \sum_{j=1}^9 X_{ji} + B_i + e_i$ ,  $e_i \sim N(0, 3)$ , and  $B_i$  is simulated as  $B_i = 0.2 \sum_{j=1}^9 X_{ji} + f_i$ ,  $f_i \sim N(0, 1)$ .

The reduced  $\mathbf{Y}|\mathbf{X}$  model is  $Y_i = 0.7 \sum_{j=1}^9 X_{ji} + B_i + k_i$ ,  $k_i \sim N(0, 4)$ . The current data sample size  $n = 200$ , replication number  $S = 10$ , and thus the synthetic data sample size  $m = nS = 2000$ .

*Results:* The results in Table 2.1 of scenario 1 show that compared to the direct MLE, the synthetic data method leads to an obvious reduction in the standard deviation of  $\gamma_X$ 's and good coverage rates of confidence intervals. In addition, it can reduce the MSE closer to the true value by 50%.

Table 2.1: Simulation results for scenario 1 with Gaussian Y, one Gaussian B and nine correlated X's. For each method, we report mean (SD) [95 % coverage rate] and MSE across 500 simulated datasets.

	Not including B	True value	Direct MLE	Synthetic Data Method
$\gamma_0$	0	0	-0.04 (0.14) [97%]	0.00 (0.18) [96%]
$\gamma_{X_1}$	0.7	0.5	0.50 (0.13) [94%]	0.51 (0.08) [95%]
$\gamma_{X_2}$	0.7	0.5	0.48 (0.13) [94%]	0.50 (0.07) [95%]
$\gamma_{X_3}$	0.7	0.5	0.49 (0.12) [95%]	0.50 (0.07) [95%]
$\gamma_{X_4}$	0.7	0.5	0.50 (0.11) [98%]	0.50 (0.07) [95%]
$\gamma_{X_5}$	0.7	0.5	0.50 (0.13) [93%]	0.50 (0.08) [93%]
$\gamma_{X_6}$	0.7	0.5	0.50 (0.12) [94%]	0.50 (0.07) [95%]
$\gamma_{X_7}$	0.7	0.5	0.50 (0.12) [94%]	0.50 (0.07) [95%]
$\gamma_{X_8}$	0.7	0.5	0.50 (0.11) [97%]	0.50 (0.07) [95%]
$\gamma_{X_9}$	0.7	0.5	0.50 (0.13) [94%]	0.51 (0.07) [95%]
$\gamma_B$	-	0.5	1.00 (0.13) [95%]	1.00 (0.11) [95%]
MSE	0.464	0.334	0.355	0.345

- **Scenario 2:** Y is binary and B is Gaussian distributed. The true model of  $Y|X, B$  is  $\text{logit}[\Pr(Y_i = 1|X_i, B_i)] = 1 + 2 \sum_{j=1}^9 X_{ji} - 3B_i$ , which gives the prevalence  $\Pr(Y = 1) \approx 0.67$ .  $B_i$  is simulated as  $B_i = 0.5 \sum_{j=1}^9 X_{ji} + e_i$ ,  $e_i \sim N(0, 0.1)$ . The current data sample size  $n = 400$ ,  $S = 10$ , and  $m = nS = 4000$ .

*Results:* As shown in Table 2.2, including B into the regression model can reduce the scaled Brier score by 15%, and improve the AUC by 9%.

Table 2.2: Simulation results for scenario 2 with binary Y, one Gaussian B and nine correlated X's. For each method, we report mean (SD) [95 % coverage rate], average scaled Brier score and AUC across 500 simulated datasets

	Not including B	True value	Direct MLE	Synthetic Data Method
$\gamma_0$	0.849	1	1.00 (0.17) [94%]	0.96 (0.08) [92%]
$\gamma_{X_1}$	0.435	2	2.01 (0.29) [96%]	1.92 (0.24) [94%]
$\gamma_{X_2}$	0.432	2	2.00 (0.28) [96%]	1.90 (0.23) [94%]
$\gamma_{X_3}$	0.437	2	2.01 (0.28) [95%]	1.90 (0.23) [95%]
$\gamma_{X_4}$	0.433	2	2.01 (0.30) [96%]	1.91 (0.24) [95%]
$\gamma_{X_5}$	0.422	2	2.02 (0.28) [95%]	1.90 (0.24) [93%]
$\gamma_{X_6}$	0.421	2	2.01 (0.28) [97%]	1.89 (0.23) [93%]
$\gamma_{X_7}$	0.431	2	2.01 (0.29) [96%]	1.91 (0.23) [94%]
$\gamma_{X_8}$	0.415	2	2.00 (0.27) [97%]	1.89 (0.23) [94%]
$\gamma_{X_9}$	0.445	2	2.01 (0.29) [96%]	1.92 (0.23) [95%]
$\gamma_B$	-	-3	-3.02 (0.45) [97%]	-2.85 (0.43) [95%]
Scaled Brier score	0.801	0.680	0.702	0.686
AUC	0.767	0.837	0.828	0.835

- **Scenario 3:** Y and B are both binary. The true model of  $Y|X, B$  is  $\text{logit}[\Pr(Y_i = 1|X_i, B_i)] = -1 + 0.2 \sum_{j=1}^4 X_{ji} - 0.2 \sum_{j=5}^7 X_{ji} - 0.5 \sum_{j=8}^9 X_{ji} + 1.5B_i$ , and  $B_i$  is simulated as  $\text{logit}[\Pr(B_i = 1|X_i)] = -0.5 + 0.5 \sum_{j=1}^5 X_{ji} - \sum_{j=6}^9 X_{ji}$ . The prevalence of Y and B are around 0.5 and 0.55, respectively, where  $n = 400$ ,  $S = 8$ , and  $m = nS = 3200$ .

*Results:* The simulation results in Table 2.3 for scenario 3, where both Y and B are binary, show that including B in the regression model can reduce the scaled Brier score by 10.8%, and increase the AUC by 5%.

Table 2.3: Simulation results for scenario 3 with binary Y, one binary B and nine correlated X's. For each method, we report mean (SD) [95 % coverage rate], average scaled Brier score and AUC across 500 simulated datasets

	Not including B	True value	Direct MLE	Synthetic Data Method
$\gamma_0$	-0.328	-1	-1.00 (0.21) [94%]	-1.00 (0.15) [96%]
$\gamma_{X_1}$	0.305	0.2	0.20 (0.13) [94%]	0.20 (0.05) [96%]
$\gamma_{X_2}$	0.318	0.2	0.21 (0.13) [95%]	0.21 (0.05) [94%]
$\gamma_{X_3}$	0.296	0.2	0.20 (0.13) [95%]	0.21 (0.05) [95%]
$\gamma_{X_4}$	0.296	0.2	0.19 (0.13) [93%]	0.20 (0.05) [95%]
$\gamma_{X_5}$	-0.066	-0.2	-0.21 (0.13) [94%]	-0.19 (0.05) [96%]
$\gamma_{X_6}$	-0.405	-0.2	-0.20 (0.13) [96%]	-0.20 (0.06) [96%]
$\gamma_{X_7}$	-0.420	-0.2	-0.19 (0.13) [96%]	-0.21 (0.06) [95%]
$\gamma_{X_8}$	-0.698	-0.5	-0.50 (0.15) [93%]	-0.51 (0.06) [96%]
$\gamma_{X_9}$	-0.713	-0.5	-0.51 (0.14) [95%]	-0.52 (0.06) [94%]
$\gamma_B$	-	1.5	1.50 (0.28) [96%]	1.49 (0.28) [95%]
Scaled Brier score	0.750	0.666	0.687	0.669
AUC	0.789	0.833	0.823	0.831

- Scenario 4:** Y is binary and two mixed types of B are included, one binary and another Gaussian. The true model of  $Y|\mathbf{X}, B_1, B_2$  is  $\text{logit}[\Pr(Y_i = 1|\mathbf{X}_i, B_{1i}, B_{2i})] = -0.2 - 0.2X_1 + 0.2 \sum_{j=2}^3 X_{ji} + 0.1 \sum_{j=4}^5 X_{ji} - 0.1X_6 - 0.3X_7 + 0.3 \sum_{j=8}^9 X_{ji} + 2B_{1i} - B_{2i}$ , of which the prevalence  $P(Y=1)$  is approximately 0.53. The binary  $B_{1i}$  is simulated as  $\text{logit}[\Pr(B_{1i} = 1|\mathbf{X}_i)] = -0.5 + 0.5 \sum_{j=1}^5 X_{ji} - \sum_{j=6}^9 X_{ji}$ , which gives the prevalence  $\Pr(B_1 = 1) \approx 0.56$ . The Gaussian  $B_{2i}$  is simulated as  $B_{2i} = 0.3 \sum_{j=1}^2 X_{ji} - 0.2 \sum_{j=3}^4 X_{ji} + 0.5 \sum_{j=5}^7 X_{ji} - 0.5 \sum_{j=8}^9 X_{ji} + e_i$ ,  $e_i \sim N(0, 0.1)$ . In this scenario,  $n = 400$ ,  $S = 8$ , and  $m = nS = 3200$ .

*Results:* The results in Table 2.4 show that the synthetic data method improves the scaled Brier score and AUC compared to the MLE, where these metrics almost attain the best possible values listed in the true value column. In addition, the coverage rates of the confidence intervals for the  $\gamma$ 's close to 95% as desired.

Table 2.4: Simulation results for scenario 4 with binary Y, binary B1 and continuous B2 and nine correlated X's. For each method, we report mean (SD) [95 % coverage rate], average scaled Brier score and AUC across 500 simulated datasets

	Not including B	True value	Direct MLE	Synthetic Data Method
$\gamma_0$	0.250	-0.2	-0.197 (0.21) [96%]	-0.13 (0.15) [93%]
$\gamma_{X_1}$	0.441	-0.2	-0.19 (0.19) [96%]	-0.14 (0.13) [94%]
$\gamma_{X_2}$	0.779	0.2	0.21 (0.21) [94%]	0.23 (0.13) [94%]
$\gamma_{X_3}$	-0.132	0.2	0.19 (0.17) [95%]	0.17 (0.10) [94%]
$\gamma_{X_4}$	-0.218	0.1	0.09 (0.17) [95%]	0.08 (0.10) [95%]
$\gamma_{X_5}$	1.047	0.1	0.10 (0.26) [96%]	0.15 (0.21) [95%]
$\gamma_{X_6}$	0.705	-0.1	-0.10 (0.28) [94%]	-0.06 (0.21) [94%]
$\gamma_{X_7}$	0.529	-0.3	-0.29 (0.27) [94%]	-0.25 (0.21) [93%]
$\gamma_{X_8}$	-0.750	0.3	0.29 (0.25) [97%]	0.23 (0.21) [95%]
$\gamma_{X_9}$	-0.757	0.3	0.30 (0.25) [96%]	0.22 (0.21) [94%]
$\gamma_{B1}$	-	1	0.996 (0.31) [96%]	0.88 (0.30) [93%]
$\gamma_{B2}$	-	2	1.99 (0.45) [95%]	1.83 (0.43) [93%]
Scaled Brier score	0.637	0.575	0.598	0.582
AUC	0.849	0.876	0.868	0.873

In summary, the results of simulation studies show that: (1) the synthetic data method can improve the efficiency of estimating  $\gamma_X$ 's while reduce the MSE of the predictions and increase the AUC for binary Y; (2) In scenario 1 where the imputation model  $B|X, Y$  is correctly specified, the estimates of  $\gamma_X$ 's and  $\gamma_B$  are unbiased; (3) In scenarios 2, 3 and 4 where the imputation model  $B|X, Y$  is misspecified, despite the improved predictive performance, there is some bias when estimating  $\gamma_X$ 's and  $\gamma_B$ . In future work, we will investigate if we can achieve even further improvements in performance using alternative or more flexible or more nonparametric approaches for imputing B.

## 2.4 Prostate Cancer Prevention Trial Data Example

To assess the performance of the synthetic data method in a real data example, we apply it to the Prostate Cancer Prevention Trial calculator. The high-grade prostate cancer calculator (PCPThg) [Thompson et al., 2006], predicts the probability of high-grade prostate cancer derived from a logistic regression based on standard clinical variables – PSA level, age, DRE findings, prior biopsy result and ethnicity. The equation for the model is:

$$\text{logit}(p_i) = -6.25 + 0.03\text{age}_i + 0.96\text{race}_i + 1.29\log(\text{PSA}_i) + 1.00\text{DRE}_i - 0.36\text{biopsy}_i. \quad (2.2)$$

where  $p_i$  is the probability of observing high-grade prostate cancer for subject  $i$  given covariates. A detailed description of the calculator and the external and internal and a validation dataset are given in Tomlins et al. [2015] and Cheng et al. [2018]. We consider incorporating two additional biomarkers that have been shown to be predictive of prostate cancer into model 2.2. One is prostate cancer antigen 3 (PCA3), a continuous variable, and the other is the indicator variable of TMPRSS2:ERG (T2:ERG) gene fusions. We consider 3 different expanded models, one with the addition of PCA3 only, one with the addition of T2:ERG only and one with the addition of both PCA3 and T2:ERG. To compare the coefficient estimation across methods, we show the estimated coefficients and standard errors in Table 2.5 from 679 observations in the internal dataset. To compare prediction power, we calculate the scaled Brier Score and the AUC based on an independent validation dataset with 1,218 observations.

Table 2.5: Results of the expanded PCPThg model. For each method, point estimate (standard error) from the internal dataset, and the scaled Brier score and the AUC from the validation dataset. The sample size of the internal dataset is 679. The sample size of the validation dataset is 1,218. Replicate number  $S=10$  giving  $m=6,790$  in the synthetic data method.

Model	PSA	Age	DRE findings	Prior biopsy history	Race	PCA3	T2:ERG	Scaled Brier Score	AUC
Original PCPThg	1.29	0.031	1.00	-0.36	0.96	–	–	0.933	0.707
Estimated PCPThg	1.06 (0.18)	0.033 (0.012)	1.15 (0.26)	-1.44 (0.27)	0.44 (0.29)	–	–	0.975	0.716
<b>Expanded model with PCA3 score</b>									
Direct regression*	0.97 (0.19)	0.009 (0.013)	1.06 (0.27)	-1.27 (0.27)	0.05 (0.31)	0.56 (0.08)	–	0.953	0.767
Synthetic data method	1.30 (0.08)	0.012 (0.006)	0.91 (0.13)	-0.56 (0.12)	0.50 (0.14)	0.57 (0.08)	–	0.878	0.765
CSPML	1.22 (0.08)	0.007 (0.005)	0.86 (0.10)	-0.20 (0.08)	0.58 (0.11)	0.56 (0.097)	–	0.888	0.759
<b>Expanded model with binary T2:ERG</b>									
Direct regression*	0.98 (0.18)	0.032 (0.012)	1.02 (0.26)	-1.41 (0.27)	0.57 (0.29)	–	0.76 (0.20)	0.930	0.744
Synthetic data method	1.21 (0.07)	0.030 (0.005)	0.96 (0.10)	-0.59 (0.09)	0.99 (0.11)	–	0.76 (0.22)	0.897	0.741
CSPML	1.14 (0.07)	0.032 (0.004)	1.06 (0.14)	-0.52 (0.11)	0.80 (0.17)	–	0.72 (0.20)	0.931	0.742
<b>Expanded model with PCA3 score and binary T2:ERG</b>									
Direct regression*	0.94 (0.19)	0.010 (0.010)	1.00 (0.28)	-1.27 (0.28)	0.15 (0.31)	0.52 (0.08)	0.47 (0.21)	0.928	0.776
Synthetic data method	1.23 (0.09)	0.008 (0.007)	0.83 (0.13)	-0.53 (0.11)	0.63 (0.15)	0.55 (0.10)	0.45 (0.20)	0.867	0.773
CSPML	1.20 (0.08)	0.008 (0.005)	0.78 (0.11)	-0.21 (0.09)	0.67 (0.12)	0.52 (0.10)	0.48 (0.27)	0.879	0.769

\*Firth corrected MLE is used

As shown in Table 2.5, for both of the expanded PCPThg models incorporating PCA3 score or binary T2:ERG, if we compare the standard errors across different methods, it is easily seen that the synthetic data method can reduce the standard errors of regression coefficients compared to direct regression by at least 50%.

The expanded PCPThg model incorporating both PCA3 score and binary T2:ERG fitted to the training dataset again shows that the proposed method can reduce the standard errors of regression coefficients compared to direct regression. The results in Table 2.5 show no improvement in AUC from using the synthetic data approach compared to direct MLE, but noticeable improvement in

the scaled Brier score.

We also include in Table 2.5 the estimates from applying the CSPML, a published method that can be applied in this case. We see it gives similar predictive performance as the synthetic data method, but the estimated coefficients differ.

## 2.5 Algebraic Justification in Two Special Cases

### 2.5.1 Estimation and Variance of $\gamma$

To establish that the synthetic data approach is asymptotically as efficient as CML approaches, we consider two special cases where closed-form results of MLE for the combined dataset of size  $n+m$  in the synthetic data approach are available, so multiple imputation does not need to be used. For these cases, we compare the synthetic data approach to the basic CML method [Cheng et al., 2018] and the CSPML method [Chatterjee et al., 2016]. These two maximum likelihood approaches are optimal based on their assumptions. The standard maximal likelihood approach based on just the observed data without incorporating external information is also provided for reference comparison. For each approach, we derive the explicit formulas for the asymptotic variance of estimated coefficients, namely,  $\hat{\gamma}$  in model (2.1).

The rationale for studying these two examples in depth is to establish some theoretical underpinning for the synthetic data approach. Given its broad applicability to other more general situations with a mixed set of continuous and categorical multivariable predictors in  $\mathbf{X}$  and  $B$ , a justification in simpler cases that can be studied analytically makes the approach more plausible in other situations where studying the analytical properties is not immediate.

As needed, we will be considering three different likelihoods, one based on the distribution  $Y|\mathbf{X}, B$ , one based on the distribution  $(Y, B)|\mathbf{X}$  and one based on the joint distribution of  $Y, \mathbf{X}$  and  $B$ . When writing distributions, we will include the parameters when necessary, e.g.  $f(Y|\mathbf{X}, B; \gamma)$ , but parameters will be excluded when not necessary.

For this study we either know the full form of the distribution of  $Y|\mathbf{X}$ , the mean of which may be characterized by a linear combination of  $\mathbf{X}$ 's, as given in equation 2.3, with known  $\beta$ 's and a known link function  $g_1$ :

$$g_1[E(Y|\mathbf{X})] = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p. \quad (2.3)$$

or we just know the mean structure but not the full distribution of  $Y|\mathbf{X}$ .

As mentioned earlier, our interest lies in building the mean structure of the  $Y|\mathbf{X}, B$  distribution as given in model 2.1. For some approaches, we will also need to consider the relationship between

$\mathbf{X}$  and  $\mathbf{B}$  for which we specify a model, the mean of which is given by

$$g_2[E(\mathbf{B}|\mathbf{X})] = \theta_0 + \theta_1 X_1 + \dots + \theta_p X_p. \quad (2.4)$$

We note that for all of these models there may be additional parameters necessary to define the full distributions (e.g. the variance  $\sigma_\beta^2$  for Gaussian  $\mathbf{Y}$ ). But for ease of notation we will not include these additional parameters unless it is necessary, thus we denote the distributions as  $f(\mathbf{Y}|\mathbf{X}; \boldsymbol{\beta})$ ,  $f(\mathbf{Y}|\mathbf{X}, \mathbf{B}; \boldsymbol{\gamma})$  and  $f(\mathbf{B}|\mathbf{X}; \boldsymbol{\theta})$ .

As a reference comparison, we will also present results for standard MLE on a complete dataset of size  $n$ . In this approach, we estimate the parameters of model 2.1 using the internal dataset of  $\mathbf{Y}$ ,  $\mathbf{X}$ , and  $\mathbf{B}$  without taking the external summary-level information into account. We obtain the estimates by maximizing the likelihood  $\prod_{i=1}^n f(\mathbf{Y}_i|\mathbf{X}_i, \mathbf{B}_i; \boldsymbol{\gamma})$  over  $\boldsymbol{\gamma}$ . Then the asymptotic covariance matrix of  $\hat{\boldsymbol{\gamma}}$  is obtained from the inverse of the Fisher information matrix.

- *Approach 1: The synthetic data method.* In special cases in which a direct solution is possible, the likelihood for a dataset of size  $n+m$  is  $\prod_{i=1}^n f(\mathbf{Y}_i, \mathbf{B}_i|\mathbf{X}_i) \prod_{i=n+1}^{n+m} f(\mathbf{Y}_i|\mathbf{X}_i)$ , and can also be written as  $\prod_{i=1}^n f(\mathbf{B}_i|\mathbf{Y}_i, \mathbf{X}_i) \prod_{i=1}^{n+m} f(\mathbf{Y}_i|\mathbf{X}_i)$ . This likelihood is then maximized over  $\boldsymbol{\gamma}$  and  $\boldsymbol{\theta}$  to obtain the MLE and the asymptotic variance is obtained from the inverse of the Fisher information.
- *Approach 2: CML on a complete dataset of size  $n$ .* For this approach we posit a model  $f(\mathbf{B}|\mathbf{X}; \boldsymbol{\theta})$  then maximize the likelihood  $\prod_{i=1}^n f(\mathbf{Y}_i, \mathbf{B}_i|\mathbf{X}_i)$ , which can be written as

$$\prod_{i=1}^n f(\mathbf{Y}_i|\mathbf{X}_i, \mathbf{B}_i; \boldsymbol{\gamma}) f(\mathbf{B}_i|\mathbf{X}_i; \boldsymbol{\theta})$$

subject to a constraint on the parameters that is derived from the external information. The equality  $f(\mathbf{Y}|\mathbf{X}; \boldsymbol{\beta}) = \int f(\mathbf{Y}|\mathbf{X}, \mathbf{B}; \boldsymbol{\gamma}) f(\mathbf{B}|\mathbf{X}; \boldsymbol{\theta}) d\mathbf{B}$  gives a relationship between the unknown parameters  $\boldsymbol{\gamma}$ ,  $\boldsymbol{\theta}$  and the known parameter  $\boldsymbol{\beta}$ . Assuming  $\boldsymbol{\theta}$  can be written as a function of  $\boldsymbol{\gamma}$  and  $\boldsymbol{\beta}$ , i.e. as  $\boldsymbol{\theta}(\boldsymbol{\gamma}, \boldsymbol{\beta})$ , then since  $\boldsymbol{\beta}$  is known the optimization problem becomes an unconstrained optimization problem, specifically maximization of

$$\prod_{i=1}^n f(\mathbf{Y}_i|\mathbf{X}_i, \mathbf{B}_i; \boldsymbol{\gamma}) f(\mathbf{B}_i|\mathbf{X}_i; \boldsymbol{\theta}(\boldsymbol{\gamma}, \boldsymbol{\beta} = \boldsymbol{\beta}^*))$$

with respect to  $\boldsymbol{\gamma}$  using the known value  $\boldsymbol{\beta}^*$  of  $\boldsymbol{\beta}$ . We consider two variations of the CML method: Approach 2.1 where only the coefficients  $\boldsymbol{\beta}$  are known, and Approach 2.2 where



both  $\beta$  and its variance  $\sigma_\beta^2$  are known.

- *Approach 3: CSPML method applied to a dataset of size  $n$ .* For this method, the estimates are obtained by maximizing the likelihood  $\prod_{i=1}^n f(Y_i, \mathbf{X}_i, B_i)$  over  $\gamma$  and the empirical distribution of  $(\mathbf{X}, B)$ , subject to a constraint [Chatterjee et al., 2016]. In this approach the distribution of  $(\mathbf{X}, B)$  is treated nonparametrically, and the constraint is derived from the integrated score equation of model 2.3. In this case the constraint is  $E_{\mathbf{X}B}[E_{Y|\mathbf{X}B}[\frac{\partial}{\partial \beta} \log\{f(Y|\mathbf{X}; \beta)\}]] = 0$ . The constrained optimization problem is implemented via Lagrange multipliers and gives both an estimate of  $\gamma$  and the non-parametric MLE of the distribution of  $(\mathbf{X}, B)$ . The asymptotic variance of  $\hat{\gamma}$  for this approach is given by  $(\mathbf{I} + \mathbf{CL}^{-1}\mathbf{C}^T)^{-1}$ , where

$$\begin{aligned}\mathbf{I} &= E_{\mathbf{X}B}[E_{Y|\mathbf{X}B}[-\frac{\partial^2}{\partial \gamma^2} \log\{f(Y|\mathbf{X}, B; \gamma)\}]] \\ \mathbf{C} &= E_{\mathbf{X}B}[E_{Y|\mathbf{X}B}[\frac{\partial}{\partial \gamma} \log\{f(Y|\mathbf{X}, B; \gamma)\} \frac{\partial}{\partial \beta} \log\{f(Y|\mathbf{X}; \beta)\}]] \\ \mathbf{L} &= E_{\mathbf{X}B}[u_\gamma(\mathbf{X}, B)u_\gamma^T(\mathbf{X}, B)]\end{aligned}$$

with  $u_\gamma(\mathbf{X}, B) = E_{Y|\mathbf{X}B}[\frac{\partial}{\partial \beta} \log\{f(Y|\mathbf{X}; \beta)\}]$ .

Intuitively, the asymptotic variance of this estimator is the inverse of information matrix  $\mathbf{I}$  of  $f(Y|\mathbf{X}, B; \gamma)$  plus the additional information due to knowing  $\beta$  from the external study  $\mathbf{CL}^{-1}\mathbf{C}^T$ .

## 2.5.2 Description of Two Special Cases

In the following two special cases, the goal is to derive the asymptotic efficiency of  $\hat{\gamma} = (\hat{\gamma}_X, \hat{\gamma}_B)^T$  through a closed-form expression for  $\text{Var}(\hat{\gamma})$ , and then compare the efficiency gain among all three approaches through the Asymptotic Relative Efficiency (ARE) of  $\text{Var}(\hat{\gamma})$ 's, compared to  $\text{Var}(\hat{\gamma})$  from the standard MLE. This will show how much efficiency we can gain by incorporating the external information from the external  $Y|\mathbf{X}$  model and what determines that gain. We describe the settings in this section and summarize the results in the subsequent section. For special case 1, the detailed algebraic derivation for each of the three approaches can be found in Appendix A.1.2, A.1.3 and A.1.4, respectively. Similarly for special case 2, the detailed algebraic derivation can be found in Appendix A.2.1, A.2.2 and A.2.3.

### 2.5.2.1 Special Case 1: Y and B are Gaussian Distributed

In this section, we assume that Y and B are continuous and have a Gaussian distribution, and assume the identity link for  $g_1$  and  $g_2$  in models 2.3 and 2.4.

Without loss of generality, we consider a simplified situation where there is only one X, i.e.  $p = 1$ . We also assume zero marginal means for Y, X and B, thus we will have no-intercept models as shown below. Let  $\sigma_X^2$  denote the variance of X. Then

$$Y|X \sim N(\beta X, \sigma_\beta^2) \quad (2.5)$$

$$Y|X, B \sim N(\gamma_X X + \gamma_B B, \sigma_\gamma^2) \quad (2.6)$$

$$B|X \sim N(\theta X, \sigma_\theta^2) \quad (2.7)$$

Depending on the information available from the external model  $Y|X$ , we consider two possible situations which correspond to two different constraints. The first situation is when the estimated coefficient  $\beta = \beta^*$  is known from model 2.5. This gives us the constraint  $\theta = \theta^* = \frac{\beta^* - \gamma_X}{\gamma_B}$ . The second situation is when both of the estimated coefficient  $\beta = \beta^*$  and the variance  $\sigma_\beta^2 = \sigma_\beta^{*2} = \gamma_B^2 \sigma_\theta^{*2} + \sigma_\gamma^2$  are known.

For the standard MLE of the complete dataset of size  $n$ , it is easy to show that the asymptotic variance of  $\hat{\gamma}_X$  and  $\hat{\gamma}_B$  are equal to  $\frac{\sigma_\gamma^2}{n\sigma_\theta^2}(\theta^2 + \frac{\sigma_\theta^2}{\sigma_X^2})$  and  $\frac{\sigma_\gamma^2}{n\sigma_\theta^2}$ , respectively.

### 2.5.2.2 Special Case 2: Y, X, B are All Binary

Assume we are interested in a saturated model:

$$\text{logit}[\Pr(Y = 1|X, B)] = \gamma_0 + \gamma_X X + \gamma_B B + \gamma_{XB} XB \quad (2.8)$$

describing the joint effect of X, B on Y, when Y, X, B are all binary variables. The external information from model 2.3 can be expressed as:

$$\text{logit}[\Pr(Y = 1|X)] = \beta_0 + \beta_1 X \quad (2.9)$$

The association between B and X is defined through the model:

$$\text{logit}[\Pr(B = 1|X)] = \theta_0 + \theta_1 X$$

We denote  $P(X = a, Y = b)$  as the probability of  $(X = a, Y = b)$  combination and  $P(B = 0, X = a, Y = b)$  as the probability of  $(B = 0, X = a, Y = b)$  combination, where  $a, b \in \{0, 1\}$ .

### 2.5.3 Summary

Based on the derivation listed in Appendix A.1 and A.2, we summarize the methods and the assumed forms of the summary-level external information for each approach in Tables 2.6 and 2.7 for special case 1 and 2, respectively.

Table 2.6: Summary of the special case 1 (Y and B are Gaussian)

Approach	Method for including external information	Available form of the external information	ARE( $\hat{\gamma}$ )*	
			$\hat{\gamma}_X$	$\hat{\gamma}_B$
Standard MLE (ref)	None	NA	1	1
1: Synthetic data method	m additional synthetic data observations	Ability to draw Y values from Y X distribution, regardless of the form	$1 - A^\dagger - \frac{\sigma_X^2 \theta^{*2}}{\sigma_X^2 \theta^{*2} + \sigma_\theta^2} D^\ddagger$	$1 - D$
2: Constrained MLE Cheng et al. [2018]	Constraint	2.1: The estimated coefficient $\beta$ is known. 2.2: Both of the estimated coefficient $\beta$ and the standard deviation $\sigma_\beta$ are known.	$1 - A$ $1 - A - \frac{\sigma_X^2 \theta^{*2}}{\sigma_X^2 \theta^{*2} + \sigma_\theta^2} D$	$1$ $1 - D$
3: CSPML Chatterjee et al. [2016]	Constraint	Known expectation of Y X	$1 - A$	1

\* ARE( $\hat{\gamma}$ ) = Var<sub>M</sub>( $\hat{\gamma}$ )/Var<sub>MLE</sub>( $\hat{\gamma}$ ), M ∈ {Synthetic Data, CML, CSPML}  
 $\dagger A = \frac{\sigma_\theta^2}{\sigma_\theta^2 + \sigma_X^2 \theta^{*2}} \frac{\sigma_\gamma^2}{\sigma_\gamma^2 + \gamma_B^2 \sigma_\theta^2}$ , where  $\theta^* = \frac{\beta^* - \gamma_X}{\gamma_B}$ .  
 $\ddagger D = \frac{2\sigma_\gamma^2 \gamma_B^2 \sigma_\theta^{*2}}{\sigma_\theta^4}$ , where  $\sigma_\theta^{*2} = \frac{\sigma_\beta^2 - \sigma_\gamma^2}{\gamma_B^2}$ .

Table 2.7: Summary of the special case 2 (Y, X and B are binary)

Approach	Method for including external information	Available form of the external information	ARE( $\hat{\gamma}$ )*		
			$\hat{\gamma}_0$	$\hat{\gamma}_X$	$\hat{\gamma}_B$ ( $\hat{\gamma}_{XB}$ )
Standard MLE (ref)	None	NA	1	1	1
1: Synthetic data method	m additional synthetic data observations	Ability to draw Y values from Y X distribution, regardless of the form	$1 - F^\dagger$	$1 - G^\ddagger$	1
2: Constrained MLE Cheng et al. [2018]	Constraint	Known estimated coefficient $\beta$	$1 - F$	$1 - G$	1
3: CSPML Chatterjee et al. [2016]	Constraint	Known expectation of Y X	$1 - F$	$1 - G$	1

\* ARE = Var<sub>M</sub>( $\hat{\gamma}$ )/Var<sub>MLE</sub>( $\hat{\gamma}$ ), M ∈ {Synthetic Data, CML, CSPML}  
 $\dagger F = \frac{\sum_{(a,b) \in \{(0,1)(0,0)\}} 1/P(X=a, Y=b)}{\sum_{(a,b) \in \{(0,1)(0,0)\}} 1/P(B=0, X=a, Y=b)}$ .  $\ddagger G = \frac{\sum_{a,b \in \{(0,1)\}} 1/P(X=a, Y=b)}{\sum_{a,b \in \{(0,1)\}} 1/P(B=0, X=a, Y=b)}$ .

Let  $\text{ARE}_M(\hat{\gamma}) = \frac{\text{Var}_M(\hat{\gamma})}{\text{Var}_{\text{MLE}}(\hat{\gamma})}$  denotes the asymptotic relative efficiency under approach M relative to the standard MLE (without external information), where  $M \in \{\text{Synthetic Data}, \text{CML}(2.1), \text{CML}(2.2), \text{CSPML}\}$  for Gaussian Y, and  $M \in \{\text{Synthetic Data}, \text{CML}, \text{CSPML}\}$  for binary Y.

**Result 1.** Special case 1: Y and B are continuous and normally distributed (Table 2.6)

- $\text{ARE}_{\text{Synthetic Data}}(\hat{\gamma}_X) = \text{ARE}_{\text{CML}(2.2)}(\hat{\gamma}_X) = 1 - A - \frac{\sigma_X^2 \theta^{*2}}{\sigma_X^2 \theta^{*2} + \sigma_\theta^{*2}} D$   
 $\text{ARE}_{\text{CSPML}}(\hat{\gamma}_X) = \text{ARE}_{\text{CML}(2.1)}(\hat{\gamma}_X) = 1 - A$
- $\text{ARE}_{\text{Synthetic Data}}(\hat{\gamma}_B) = \text{ARE}_{\text{CML}(2.2)}(\hat{\gamma}_B) = 1 - D$   
 $\text{ARE}_{\text{CSPML}}(\hat{\gamma}_B) = \text{ARE}_{\text{CML}(2.1)}(\hat{\gamma}_B) = 1$

where  $A = \frac{\sigma_\theta^2}{\sigma_\theta^2 + \sigma_X^2 \theta^{*2}} \frac{\sigma_\gamma^2}{\sigma_\gamma^2 + \gamma_B^2 \sigma_\theta^2}$ ,  $D = \frac{2\sigma_\gamma^2 \gamma_B^2 \sigma_\theta^{*2}}{\sigma_\beta^{*4}}$ ,  $\theta^* = \frac{\beta^* - \gamma_X}{\gamma_B}$  and  $\sigma_\theta^{*2} = \frac{\sigma_\beta^{*2} - \sigma_\gamma^2}{\gamma_B^2}$ . In summary, the synthetic data method has the same asymptotic variance as the CML (approach 2.2), and both are more efficient than the CSPML and the CML (approach 2.1). For  $\hat{\gamma}_X$  the CSPML and the CML (approach 2.1) are more efficient than the standard MLE. For  $\hat{\gamma}_B$  the CSPML and the CML (approach 2.1) have the same efficiency as the standard MLE.

**Result 2.** Special case 2: Y, X, and B are all binary (Table 2.7)

- $\text{ARE}_{\text{Synthetic Data}}(\hat{\gamma}_0) = \text{ARE}_{\text{CML}}(\hat{\gamma}_0) = \text{ARE}_{\text{CSPML}}(\hat{\gamma}_0) = 1 - F$
- $\text{ARE}_{\text{Synthetic Data}}(\hat{\gamma}_X) = \text{ARE}_{\text{CML}}(\hat{\gamma}_X) = \text{ARE}_{\text{CSPML}}(\hat{\gamma}_X) = 1 - G$
- $\text{ARE}_{\text{Synthetic Data}}(\hat{\gamma}_B) = \text{ARE}_{\text{CML}}(\hat{\gamma}_B) = \text{ARE}_{\text{CSPML}}(\hat{\gamma}_B) = 1$
- $\text{ARE}_{\text{Synthetic Data}}(\hat{\gamma}_{XB}) = \text{ARE}_{\text{CML}}(\hat{\gamma}_{XB}) = \text{ARE}_{\text{CSPML}}(\hat{\gamma}_{XB}) = 1$

where  $F = \frac{\sum_{(a,b) \in \{(0,1)(0,0)\}} 1/P(X=a, Y=b)}{\sum_{(a,b) \in \{(0,1)(0,0)\}} 1/P(B=0, X=a, Y=b)}$  and  $G = \frac{\sum_{a,b \in \{0,1\}} 1/P(X=a, Y=b)}{\sum_{a,b \in \{0,1\}} 1/P(B=0, X=a, Y=b)}$ . In conclusion, the synthetic data method, CML and CSPML all converge to the same asymptotic variance. We notice that  $\hat{\gamma}_0$  and  $\hat{\gamma}_X$  are more efficient than the standard MLE while  $\hat{\gamma}_B$  and  $\hat{\gamma}_{XB}$  have the same efficiency as the standard MLE.

## 2.5.4 Justification from Another Perspective

In the two special cases, we show that using the synthetic data approach with very large m (i.e. total size of synthetic data) gives identical asymptotic variance for the parameters of model 2.1 as the constrained ML approach. Below we provide a different intuitive justification for the synthetic data approach, for a more general situation, if certain conditions apply. Assume that Y

and  $B$  are scalar random variables and that  $\mathbf{X}$  is a vector of covariates. We will assume parametric models for all the conditional distributions, and these can be written as  $f(Y, B|\mathbf{X}, \phi)$ ,  $f(Y|\mathbf{X}, B; \gamma)$ ,  $f(Y|\mathbf{X}; \beta)$ ,  $f(B|\mathbf{X}; \theta)$  and  $f(B|\mathbf{X}, Y; \kappa)$ . Assume that  $f(Y|\mathbf{X}, B; \gamma)$  is the model of interest, and that  $f(Y|\mathbf{X}; \beta)$  is the form of the model that was fit to the external data, and that the estimate of  $\beta$  from the external data approximates the true value of  $\beta$ . We assume that all these models represent the true distributions and are compatible with each other in the sense that  $f(Y, B|\mathbf{X}, \phi) = f(Y|\mathbf{X}, B; \gamma) \times f(B|\mathbf{X}; \theta) = f(B|\mathbf{X}, Y; \kappa) \times f(Y|\mathbf{X}; \beta)$ . We assume there is a 1-to-1 mapping between  $\phi$  and  $(\gamma, \theta)$  and between  $\phi$  and  $(\kappa, \beta)$ , and that  $\kappa$  and  $\beta$  are distinct and that  $\gamma$  and  $\theta$  are distinct. With these conditions, we can write  $f(Y, B|\mathbf{X}; \phi)$  as  $f(Y, B|\mathbf{X}; \kappa, \beta)$ .

With this set-up, the CML estimate is obtained by maximizing the likelihood  $\prod_{i=1}^n f(Y_i, B_i|\mathbf{X}_i; \phi)$  over  $\phi$ , subject to the known  $\beta$ . This can be rewritten as maximizing the likelihood  $\prod_{i=1}^n f(Y_i, B_i|\mathbf{X}_i; \kappa, \beta)$  over  $\kappa$ , subject to the known  $\beta$ . Then from the combination of the estimate of  $\kappa$  and the known  $\beta$  we can obtain the estimate of  $\gamma$ .

The synthetic data method consists of maximizing the likelihood

$$\prod_{i=1}^n f(Y_i, B_i|\mathbf{X}_i; \phi) \prod_{i=n+1}^{n+m} f(Y_i|\mathbf{X}_i; \beta)$$

which is equivalent to maximizing

$$\prod_{i=1}^n f(Y_i, B_i|\mathbf{X}_i; \kappa, \beta) \prod_{i=n+1}^{n+m} f(Y_i|\mathbf{X}_i; \beta)$$

over  $\kappa$  and  $\beta$ . When optimizing over  $\beta$  for fixed  $\kappa$ , the second term  $\prod_{i=n+1}^{n+m} f(Y_i|\mathbf{X}_i; \beta)$  will dominate the optimization procedure when  $m$  is very large. Thus the estimate of  $\beta$  will essentially reproduce the known value from the external data (since this was the value used to generate the synthetic data). Thus the synthetic data method will reduce to the maximization of the remaining part of the likelihood  $\prod_{i=1}^n f(Y_i, B_i|\mathbf{X}_i; \kappa, \beta)$  with  $\beta$  fixed, which is exactly the CML method.

The requirement that all the conditional distributions are compatible with each other will not usually be true, but maybe a reasonable approximation if flexible enough models are being used. The conditions do hold for the normal and the tri-binary examples in Section 2.5.2.1 and 2.5.2.2. Another case where they hold is when  $Y$  and  $B$  follow a bivariate normal distribution given  $\mathbf{X}$ , i.e.  $Y, B|\mathbf{X} \sim N \left( \begin{pmatrix} \beta\mathbf{X} \\ \theta\mathbf{X} \end{pmatrix}, \begin{pmatrix} \sigma_\beta^2 & \rho\sigma_\beta\sigma_\theta \\ \rho\sigma_\beta\sigma_\theta & \sigma_\theta^2 \end{pmatrix} \right)$ . Then the CML is to maximize the likelihood

$\prod_{i=1}^n f(Y_i, B_i|\mathbf{X}_i; \beta, \theta, \sigma_\beta, \sigma_\theta, \rho)$  over  $\theta, \sigma_\theta$  and  $\rho$  subject to the known  $\beta$  and  $\sigma_\beta$ .

## 2.6 Discussion

In this chapter, we introduced the synthetic data method for incorporating summary-level information from well-established external models into the regression model estimation based on internal data. We demonstrate in some special cases that with a large number of synthetic data observations, the synthetic data approach is asymptotically as efficient as the CML approach. This provides some justification for what at first sight might seem to be an ad-hoc approach. In a simulation study, we demonstrate the ability of the method to improve the predictive ability of the model.

A key advantage of the synthetic data method is that it naturally incorporates the prior knowledge into the internal data by creating large “fake” data that is compatible with the  $Y|X$  established model. By creating pseudo-data from  $Y|X$  instead of constrained optimization, the synthetic data method not only simplified the task from solving complex constrained optimization, but also provides a potentially more flexible and general framework to handle this problem. The only requirement for the synthetic data approach is the ability to generate  $Y$  values given  $X$  from the information of the external models, without the need to know the exact form of the model. It is broadly applicable for general data types for  $Y$ ,  $X$  and  $B$ , and when  $B$  is more than one new biomarker. It can be extended to the situation where more than one external model is available, i.e.  $Y|X_1, Y|X_2, \dots, Y|X_k$ . In this setting, a combination of external studies that measured overlapping but necessarily identical covariates can provide joint information to develop a model for  $Y|X$  model, where  $X$  is the union of  $X_1, X_2, \dots, X_k$ . We will elaborate this in Chapter 4.

The CSPML approach is also broadly applicable, and can handle multiple  $B$ 's, and it has some optimality properties. But it does require knowledge of the form of the  $Y|X$  model and requires that the distribution of the  $X$ 's are identical in the external and the internal populations, which seems unlikely to be satisfied in practice. When analyzing the synthetic dataset, the value of  $B$  can be considered as missing, which converts the problem of incorporating external information into a problem of analyzing data with missing values. If multiple imputation procedures are to be used to impute the value of  $B$ , then further research would be needed to suggest efficient and robust ways in which this should be implemented. There is the potential to improve even further on the method by using different ways of imputing  $B$ , beyond the approach we illustrated in the simulation study.

Another interesting issue that will need to be investigated is the size of  $m$ . The theoretical result in this chapter suggests that  $m$  should be very large, but this is under the assumption that the  $Y|X$  and  $Y|X, B$  models are compatible with each other. In practice, they are unlikely to be exactly compatible, which would suggest limiting the size of  $m$ . A pragmatic suggestion is to make  $m$  equal to the size of the external data, if that is known. By doing this the amount of information in the synthetic data about the relationship between  $Y$  and  $X$  is similar to the amount of information in the external data about the relationship between  $Y$  and  $X$ .

## 2.7 Publication

The content of this chapter has been published in *Canadian Journal of Statistics* at DOI: 10.1002/cjs.11513.

## CHAPTER 3

# A Meta-Inference Framework to Integrate Multiple External Models into a Current Study

### 3.1 Introduction

In Chapter 2, we introduced a synthetic data method as a flexible alternative to the constrained maximum likelihood (CML) approach. In two special cases and extensive simulation studies, we showed that the synthetic data method had the same asymptotic properties as a constrained semi-parametric maximum likelihood (CSPML) approach [Chatterjee et al., 2016], who converted the external summary-level information into a constraint and then maximized the internal data likelihood subject to this constraint.

Several methods were built upon the work of Chatterjee et al. [2016], e.g. Zhang et al. [2020] extended CSPML to account for parameter uncertainty of the external model while Kundu, Tang, and Chatterjee [2019] extended CSPML to a generalized meta-analysis approach. However, the CSPML method requires the joint distribution of  $(Y, X, B)$  to be the same in the internal and the external population, a strong assumption, which although unverifiable, we expect would be frequently violated, and can cause bias when violated. Estes, Mukherjee, and Taylor [2017] later proposed a matrix-weighted average remedy by constructing an empirical Bayes (EB) estimator that can reduce the potential bias. As an extension and adaption of Estes et al. [2017], in this chapter, we propose a meta-inference framework using a composite of EB estimators to accommodate the situation where multiple external prediction models are available to help improve the inference of the current study.

We consider the situation in which there are  $K$  external studies ( $K \geq 2$ ), each of which developed a prediction model for the same outcome. The parameter estimates of the external models are known, but the individual-level data are not available. The goal is to develop a prediction model that uses all the possible covariates, using data from an internal study and the parameter estimates from the external models. The parameters of this model are the quantities of interest. Each of the external studies may use a slightly different set of covariates but the internal data are assumed



to contain all available covariates, as well as the new biomarkers that are not included in any of the external models. We propose a meta-inference framework using an empirical Bayes estimation approach, which first separately incorporates the different summary information from each external study into the internal study, and then takes a weighted average of the resulting estimators to give a final overall estimate of the parameters of interest. We show that the proposed final estimators are more efficient than the simple analysis of the internal data, as well as outperform the estimators that integrate the information from a single external model.

The rest of this chapter is organized as follows: In Section 3.2, we introduce two key existing estimators before introducing the proposed methodology along with two weighted estimators, followed by corresponding large sample results. In Section 3.3, we demonstrate the potential of our proposed method through a simulation study. In Section 3.4, we apply the proposed method to predict the risk of high-grade prostate cancer incorporating information from two existing risk prediction models. We present a discussion in Section 3.5.

## 3.2 Models and Methods

### 3.2.1 General Description of the Problem

Let  $Y$  denote the outcome of interest, which can be either continuous or binary.  $\mathbf{X}$  is a set of  $p$  standard variables and let  $B$  denote a new biomarker. Our target of interest is the mean structure of  $Y|\mathbf{X}, B$ :

$$g(E(Y|\mathbf{X}, B)) = \mathbf{X}\boldsymbol{\gamma}_X + B\gamma_B = \gamma_{X_0} + \gamma_{X_1}X_1 + \dots + \gamma_{X_p}X_p + \gamma_B B, \quad (3.1)$$

where  $g$  is the known link function. We assume that a small dataset of size  $n$  with variables  $Y, \mathbf{X}$  and a new covariate  $B$  is available to us for building model 3.2.1. For each external study  $k$ ,  $k \in \{2, \dots, K\}$ , a prediction model for the same outcome  $Y$  has been built using predictors  $\mathbf{X}_k$ , a subset of the internal  $\mathbf{X}$ . Each external model may use slightly different predictors to predict  $Y$ :

$$g(E(Y|\mathbf{X}_k)) = \mathbf{X}_k\boldsymbol{\beta}_k = \beta_0 + \beta_1X_1 + \dots + \beta_{p_k}X_{p_k},$$

where  $p_k \leq p$  is the dimension of  $\mathbf{X}_k$ . We assume that the distribution of  $Y|\mathbf{X}, B$  is correctly specified, but the external  $Y|\mathbf{X}_k$  distributions need not be.

We assume  $K$  large, well-characterized previous studies from the external populations describe the provided information on the calculated distribution of  $Y|\mathbf{X}_k$ . These information will come in the forms of estimated model parameters  $\hat{\boldsymbol{\beta}}_k$ . The goal is to develop a framework, in which we can utilize all  $K$  external  $\hat{\boldsymbol{\beta}}_k$ 's to improve the estimation efficiency of the internal study.

We introduce some of the important notation that will frequently appear in later sections:

- $f_\beta(Y|X_k)$ : the study-specific distribution of the  $k^{\text{th}}$  external model  $Y|X_k$ ;
- $f_\gamma(Y|X, B)$ : distribution of the target model  $Y|X, B$ ;
- $\hat{\gamma}_I$ : the unconstrained estimator by using the internal data only;
- $\hat{\gamma}_{\text{CSPML}}$ : the constrained semi-parametric maximum likelihood estimator (CSPML) proposed by Chatterjee et al. [2016];
- $\hat{\gamma}_{\text{EB}}$ : the empirical Bayes (EB) estimator proposed by Estes et al. [2017].

### 3.2.2 Two Existing Estimators

The proposed method was developed on the foundation of two existing methods, the CSPML approach [Chatterjee et al., 2016] and the EB approach [Estes et al., 2017]. The CSPML estimator considered the same problem described here in a special case where  $K=1$ ; and EB estimator applied the empirical Bayes method to CSPML, calibrating the potential bias due to non-transportability. Therefore, it is necessary to introduce these two core methods first before considering the  $K > 1$  situation.

In the CSPML, the proposed estimator  $\hat{\gamma}_{\text{CSPML}}$  incorporates the external regression coefficients to calibrate the current regression model. Denote  $U_\beta(Y|X)$  as the score function of the external  $Y|X; \beta$  model. It converts the external model parameter  $\hat{\beta}$  to a constraint by connecting the external score function with the target distribution  $f_\gamma(Y|X, B)$ :

$$\begin{aligned} 0 &= E_{Y,X,B}[U_\beta(Y|X)] = E_{X,B}\{E_{Y|X,B}[U_\beta(Y|X)|X, B]\} \\ &= \int_{X,B} \int_{Y|X,B} U_\beta(Y|X) f_\gamma(Y|X, B) dY dF(X, B) \\ &= \sum_{i=1}^n \int_{Y|X,B} U_\beta(Y|X) f_\gamma(Y|X, B) dY p_i, \end{aligned}$$

where  $dF(X, B)$  is the empirical probability distribution  $p_i = \Pr(X = X_i, B = B_i)$  for the internal observations and  $\sum_{i=1}^n p_i = 1$ . Then  $\hat{\gamma}_{\text{CSPML}}$  was obtained using Lagrange multipliers by solving the following Lagrangian function:

$$\hat{\gamma}_{\text{CSPML}} = \arg \max_{\gamma, p_i} \left\{ \prod_{i=1}^n f_\gamma(Y_i|X_i, B_i) p_i + \lambda_1 \left( \sum_{i=1}^n p_i - 1 \right) + \lambda_2 \sum_{i=1}^n \int U_\beta(Y|X) f_\gamma(Y|X, B) dY p_i \right\}$$

Chatterjee et al. [2016] provided the asymptotic variance of  $\hat{\gamma}_{\text{CSPML}}$ , showing the efficiency gain of  $\hat{\gamma}_{\text{CSPML}}$  compared to  $\hat{\gamma}_I$ .

Estes et al. [2017] showed that in the CSPML approach, the strict assumption of the identical joint distribution of  $(Y, X, B)$  between the internal and the external population (also known as the full transportable assumption of the joint distribution of  $[Y, X, B]$ ) is hard to satisfy in reality, and ignoring it can lead to substantial bias. Assuming the target conditional distribution  $(Y|X, B; \gamma)$  is correctly specified in the internal study and the underlying true parameter  $\gamma$  follows a stochastic framework, Estes et al. [2017] proposed an empirical Bayes (EB) estimator  $\hat{\gamma}_{\text{EB}}$ , which can be viewed as a matrix-weighted average of the internal estimate  $\hat{\gamma}_I$  and the CSPML estimator  $\hat{\gamma}_{\text{CSPML}}$ . The EB estimator uses the difference  $\hat{\gamma}_I - \hat{\gamma}_{\text{CSPML}}$  to measure the distributional similarity of the joint distribution  $(Y, X, B)$  between the internal and the external population, and will down-weight  $\hat{\gamma}_{\text{CSPML}}$  if the lack of full transportability leads to a poor estimate. Therefore, the EB estimator is robust to departures from full transportability assumption in specific external populations. The features and assumptions for the CSPML, the EB and the proposed estimators are summarized in Appendix B.1.

Specifically, the EB approach first posits a stochastic framework connecting the internal estimator  $\hat{\gamma}_I$  and the underlying true parameter  $\gamma \sim N(\gamma_0, \mathbf{A})$  for some covariance matrix  $\mathbf{A}$ . Since 
$$\begin{cases} \gamma \sim N(\gamma_0, \mathbf{A}) \\ \hat{\gamma}_I | \gamma \sim N(\gamma, \Sigma) \end{cases},$$
 the posterior Bayes estimate of  $\gamma$  equals to  $\mathbf{A}(\Sigma + \mathbf{A})^{-1}\hat{\gamma}_I + \Sigma(\Sigma + \mathbf{A})^{-1}\gamma_0$ . Replacing  $\gamma_0$  with the CSPML estimator  $\hat{\gamma}_{\text{CSPML}}$  and empirically estimating  $\mathbf{A}$  and  $\Sigma$ , we obtain the EB estimator  $\hat{\gamma}_{\text{EB}} = \hat{\mathbf{A}}(\hat{\Sigma} + \hat{\mathbf{A}})^{-1}\hat{\gamma}_I + \hat{\Sigma}(\hat{\Sigma} + \hat{\mathbf{A}})^{-1}\hat{\gamma}_{\text{CSPML}} \stackrel{\text{def}}{=} \hat{\mathbf{W}}\hat{\gamma}_I + (\mathbf{I} - \hat{\mathbf{W}})\hat{\gamma}_{\text{CSPML}}$ , where  $\hat{\mathbf{A}} = (\hat{\gamma}_I - \hat{\gamma}_{\text{CSPML}})(\hat{\gamma}_I - \hat{\gamma}_{\text{CSPML}})^T$  quantifies the difference between  $\hat{\gamma}_I$  and  $\hat{\gamma}_{\text{CSPML}}$ ,  $\hat{\Sigma}$  is the MLE of the variance of  $\hat{\gamma}_I$ , and  $\hat{\mathbf{W}} = \hat{\mathbf{A}}(\hat{\Sigma} + \hat{\mathbf{A}})^{-1}$  is the empirical weights. Therefore, the EB estimator can be viewed as a matrix generalization of a weighted average of vectors  $\hat{\gamma}_I$  and  $\hat{\gamma}_{\text{CSPML}}$ , of which the empirical weights will shrink the EB estimator towards  $\hat{\gamma}_I$  when  $\hat{\gamma}_{\text{CSPML}}$  is biased. The EB estimator's shrinkage effect limits the impact of external model information that is not compatible with the internal data and thus protects against the severe bias. Furthermore, when the joint distribution of  $(Y, X, B)$  is similar in the two populations, more precision will be gained by incorporating the external model.

### 3.2.3 Proposed Meta-Framework for Inference

We build upon the empirical Bayes work by Estes et al. [2017] and generalize it to accommodate the situation where we can combine multiple external model estimates into the internal study. Our proposed meta-inference framework allows each established model to have different dimensions. The framework consists of generating an EB estimator for each of the external  $\hat{\beta}_k$  from the fitted

regression  $Y|X_k$ , and then constructing the final estimates of the target through a weighted average, considering the correlation structure among all the EB estimators.

The proposed framework contains two steps:

- Step 1: For each of the  $K$  external model estimates  $\hat{\beta}_k$ , first apply the CSPML method [Chatterjee et al., 2016], and then apply the EB method [Estes et al., 2017]:

$$\text{Internal data} + \text{External } \hat{\beta}_k \xrightarrow{\text{CSPML}} \hat{\gamma}_{\text{CSPML}_k} \xrightarrow{\text{EB}} \hat{\gamma}_{\text{EB}_k} = \hat{W}_k \hat{\gamma}_I + (I - \hat{W}_k) \hat{\gamma}_{\text{CSPML}_k},$$

after which we obtain a total of  $K$   $\hat{\gamma}_{\text{EB}}$ 's.

- Step 2: We propose two estimators of  $\gamma$  to be the weighted average of  $\hat{\gamma}_{\text{EB}}$ 's, of the form  $\sum_{k=1}^K w_k \hat{\gamma}_{\text{EB}_k}$ , in which each element of the final estimate is the weighted average of the  $K$  separate estimates for that element. One composite estimator is called the optimal covariance weighted estimator (OCWE) and another is called the selective coefficient learner (SC-Learner).

In step 1, we separately integrate each of the external  $\hat{\beta}$ 's with the internal data, and the EB method accounts for the potential bias caused by the heterogeneity of the internal and that specific external population. This first step also unifies the disparate dimensions of the external models to be the same as the target model 3.1, and improves the efficiency of parameter estimates for those covariates that were used in the external models.

In step 2, the challenge is to combine  $K$  correlated vectors of EB estimators while maximizing the efficiency gain of the overall prediction. The simplest, yet not the most attractive, solution is to naively average  $\hat{\gamma}_{\text{EB}}$ 's, i.e.  $\frac{1}{K} \sum_{k=1}^K \hat{\gamma}_{\text{EB}_k}$ . Better weighting approaches take into account the variance and/or correlation among  $\hat{\gamma}_{\text{EB}}$ 's. One option is to use the inverse of the prediction variances as weights, i.e.  $w_k = \frac{1/\sum_{i=1}^n \hat{\text{Var}}[(X_i, B_i) \hat{\gamma}_{\text{EB}_k}]}{\sum_{k=1}^K 1/\sum_{i=1}^n \hat{\text{Var}}[(X_i, B_i) \hat{\gamma}_{\text{EB}_k}]}$ , with the same  $w_k$  used for all elements of  $\gamma$ . This method incorporates the within-estimator variance while ignoring the between-estimator covariance (i.e. ignoring the fact that the  $\hat{\gamma}_{\text{EB}}$ 's are not independent). Other popular design criteria that seek the optimal  $\hat{\gamma}$  to minimize the variance-covariance matrix of  $\gamma$  include D-optimality that minimizes the determinant of the matrix, A-optimality that minimizes the trace of the matrix, I-optimality (also known as V- or IV- or Q-optimality) that minimize the average prediction variance and G-optimality that minimizes the maximum prediction variance [Goos and Jones, 2011]. Since all criteria had similar performance in this study, we consider an adaptive version of I-optimality which seeks  $\hat{\gamma}$  that minimizes the average variance of the predicted estimator  $\frac{1}{n} \sum_{i=1}^n \hat{\text{Var}}[(X_i, B_i) \hat{\gamma}]$ . We propose two weighted estimators that accounts for both within and between variances among  $\hat{\gamma}_{\text{EB}}$ 's:

1. **The optimal covariance weighted estimator (OCWE):** OCWE views each  $\hat{\gamma}_{EB}$  as a whole and provides the same weight  $w_k$  to each covariate coefficient within  $\hat{\gamma}_{EB_k}$  that minimizes the overall estimated prediction variance, i.e.

$$\hat{\gamma}_{OCWE} = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{i=1}^n \hat{\operatorname{Var}}[(\mathbf{X}_i, B_i)\hat{\gamma}(\mathbf{w})],$$

where  $\hat{\gamma}(\mathbf{w}) = \sum_{k=1}^K w_k \hat{\gamma}_{EB_k}$  and  $\mathbf{w} = (w_1, \dots, w_K)^T$  denotes the positive weights that add up to one.

2. **The selective coefficient learner (SC-Learner):** Instead of seeking a fixed weight for each  $\hat{\gamma}_{EB}$  as in OCWE, SC-Learner attempts to find a set of weights separately for each covariate coefficient (from intercept  $\hat{\gamma}_0$  to slopes  $\hat{\gamma}_{X_1}, \dots, \hat{\gamma}_{X_p}, \hat{\gamma}_B$ ) that minimize the corresponding variance, and thus each coefficient in one  $\hat{\gamma}_{EB}$  can have different weights. Let  $E_j$  denotes the index set of the external models that included the predictor  $X_j$ . For each predictor  $X_j$  ( $j \in 0, \dots, p$ ), SC-Learner first selects  $\hat{\gamma}_{X_j^k}$  from  $\hat{\gamma}_{EB_k}$  which used  $X_j$  as a predictor in the external model, and then uses the inverse variance as the weight  $w_{kj} = \frac{1/\hat{\operatorname{Var}}(\hat{\gamma}_{X_j^k})}{\sum_{k \in E_j} 1/\hat{\operatorname{Var}}(\hat{\gamma}_{X_j^k})}$ . The final estimate of each  $\gamma_{X_j}$  is an inverse variance-weighted estimator using selective coefficients from  $\hat{\gamma}_{EB}$ 's:  $\hat{\gamma}_{X_j}^* = \sum_{k \in E_j} w_{kj} \hat{\gamma}_{X_j^k}$ . We will use  $\hat{\gamma}_B$  from the direct regression  $\hat{\gamma}_I$  as the final estimate for the  $B$  variable, since  $B$  is only available from the internal data and no external models have used  $B$  as predictors. Thus,

$$\hat{\gamma}_{SC-Learner} = [\hat{\gamma}_{X_0}^*, \hat{\gamma}_{X_1}^*, \dots, \hat{\gamma}_{X_p}^*, \hat{\gamma}_B]^T$$

To illustrate this method, consider a hypothetical example with three external models—Model 1  $Y|X_1, X_2, X_3$ , Model 2  $Y|X_1, X_2$ , and Model 3  $Y|X_1, X_3$  available to build the target model  $Y|X_1, X_2, X_3, B$  together with the internal data. When considering the final estimated coefficient of  $X_3$ ,  $\hat{\gamma}_{X_3^1}$  and  $\hat{\gamma}_{X_3^3}$  from the external models 1 and 3 will be used, while  $\hat{\gamma}_{X_3^2}$  that did not add extra information to  $X_3$  will be excluded.

In both proposed estimators, we use the asymptotic variance-covariance structure derived from the large sample theory to capture the correlation among  $\hat{\gamma}_{EB's}$ , which will be discussed in detail in Section 3.2.4.

### 3.2.4 Asymptotic Normality and Large Sample Results

The following proposition extends the asymptotic normality of the CSPML estimator [Chatterjee et al., 2016] to higher dimension, as well as showing the correlation structure between  $\hat{\gamma}_{CSPML}$  and

$\hat{\gamma}_I$ .

*Proposition 1.* Let  $\hat{\boldsymbol{\eta}} = (\hat{\gamma}_{\text{CSPML}_1}^T, \dots, \hat{\gamma}_{\text{CSPML}_K}^T, \hat{\gamma}_I^T)^T$ , and  $\boldsymbol{\eta}_0 = (\gamma_0^T, \dots, \gamma_0^T, \gamma_0^T)^T$  with  $\gamma_0$  the true value of  $\gamma_{\text{CSPML}}$  and  $\gamma_I$ . Under regularity conditions described in Chatterjee et al. [2016], as the internal sample size  $n \rightarrow \infty$ ,  $\sqrt{n}(\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}_0)$  converges in distribution to a normal distribution with zero-mean and covariance matrix given by

$$\begin{pmatrix} (B + C_1^T L_{11}^{-1} C_1)^{-1} & \dots & \dots & \dots & \dots & (B + C_1^T L_{11}^{-1} C_1)^{-1} \\ & \ddots & & \vdots & & \vdots \\ & (B + C_j^T L_{jj}^{-1} C_j)^{-1} & \dots & \text{Cov}(\hat{\gamma}_{\text{CSPML}_j}, \hat{\gamma}_{\text{CSPML}_k}) & \dots & (B + C_j^T L_{jj}^{-1} C_j)^{-1} \\ & & \ddots & \vdots & & \vdots \\ & & & (B + C_K^T L_{KK}^{-1} C_K)^{-1} & \dots & (B + C_K^T L_{KK}^{-1} C_K)^{-1} \\ & & & & & B^{-1} \end{pmatrix}, \quad (3.2)$$

where  $B = -E\{\frac{\partial^2}{\partial \gamma^T \partial \gamma} \log f_\gamma(Y|X, B)\}$ ,  $u_\gamma(\beta_j) = E_{Y|X, B}\{\frac{\partial}{\partial \beta_j} \log f_{\beta_j}(Y|X)\}$ ,  $C_j = E\{\frac{\partial}{\partial \gamma} u_\gamma(\beta_j)\}$ ,  $L_{jk} = E\{u_\gamma^T(\beta_j) u_\gamma(\beta_k)\}$ ,  $\text{Cov}(\hat{\gamma}_{\text{CSPML}_j}, \hat{\gamma}_{\text{CSPML}_k}) = (B + C_j^T L_{jj}^{-1} C_j)^{-1} (B + C_j L_{jj}^{-1} L_{jk} L_{kk}^{-1} C_k) (B + C_k^T L_{kk}^{-1} C_k)^{-1}$ , and  $j, k \in \{1, \dots, K\}$ .

As shown in *Proposition 1*, as well as in Chatterjee et al. [2016], the asymptotic variance of  $\hat{\gamma}_{\text{CSPML}}$  will decrease from  $B^{-1}$  to  $(B + C_k^T L_{kk}^{-1} C_k)^{-1}$  after incorporating the  $k^{\text{th}}$  external model information. We further show additional two conclusions: (i) the covariance between  $\hat{\gamma}_{\text{CSPML}_k}$  and  $\hat{\gamma}_I$  is equivalent to the variance of  $\hat{\gamma}_{\text{CSPML}_k}$ , i.e.  $\text{Cov}(\hat{\gamma}_{\text{CSPML}_k}, \hat{\gamma}_I) = \text{Var}(\hat{\gamma}_{\text{CSPML}_k})$ ; (ii) the covariance between any two  $\hat{\gamma}_{\text{CSPML}_j}$  and  $\hat{\gamma}_{\text{CSPML}_k}$  equals to  $(B + C_j L_{jj}^{-1} L_{jk} L_{kk}^{-1} C_k)$  multiplied by their variances, i.e.  $\text{Cov}(\hat{\gamma}_{\text{CSPML}_j}, \hat{\gamma}_{\text{CSPML}_k}) = (B + C_j^T L_{jj}^{-1} C_j)^{-1} (B + C_j L_{jj}^{-1} L_{jk} L_{kk}^{-1} C_k) (B + C_k^T L_{kk}^{-1} C_k)^{-1}$ . In Appendix B.2, we show the extension of *Proposition 1* when the uncertainty of the external  $\hat{\beta}$  is known. As expected, this modification makes a difference only when the uncertainty is large.

*Proposition 2.* Let  $Z \equiv \hat{\gamma}_I - \hat{\gamma}_{\text{CSPML}}$  and  $\hat{V}_I \equiv \text{Var}(\hat{\gamma}_I)$ . We can re-parameterize  $\hat{\gamma}_{\text{EB}}$  as a function of  $Z$  and  $\hat{\gamma}_{\text{CSPML}}$ :

$$\hat{\gamma}_{\text{EB}} = \hat{\gamma}_{\text{CSPML}} + Z \left( 1 - \frac{1}{1 + Z^T \hat{V}_I^{-1} Z} \right),$$

where  $Z^T \hat{V}_I^{-1} Z$  is a scalar. Equivalently,  $\hat{\gamma}_{\text{EB}}$  can be written as a function of  $Z$  and  $\hat{\gamma}_I$ , i.e.  $\hat{\gamma}_{\text{EB}} = \hat{\gamma}_I - Z \frac{1}{1 + Z^T \hat{V}_I^{-1} Z}$ . The proof is listed in Appendix B.3.

In the proposed method, we use the asymptotic variance-covariance structure derived from *Proposition 1* and *2* to capture the correlation among  $\hat{\gamma}_{\text{EB's}}$ . Under the assumption that the external population is representative of the target population of interest, the mean of  $Z$  converges to zero. In Appendix B.3, we show that  $Z$  and  $\hat{\gamma}_{\text{CSPML}}$  are independently normal-distributed with closed-form mean and variance, from which it is easy to simulate many values

of  $Z$  and  $\hat{\gamma}_{\text{CSPML}}$ . Therefore, according to *Proposition 2* and denote  $f(Z)$  as  $Z(1 - \frac{1}{1+Z^T\hat{V}_I^{-1}Z})$ , we can easily obtain the numeric values of  $\text{Var}(\hat{\gamma}_{\text{EB}_k})$  through the equation  $\text{Var}(\hat{\gamma}_{\text{EB}_k}) = \text{Var}(\hat{\gamma}_{\text{CSPML}_k}) + \text{Var}[f(Z)] + 2\text{Cov}[\hat{\gamma}_{\text{CSPML}_k}, f(Z)]$ , using the simulated values of  $Z$  and  $\hat{\gamma}_{\text{CSPML}}$ . A similar idea can be applied to obtain  $\text{Cov}(\hat{\gamma}_{\text{EB}_j}, \hat{\gamma}_{\text{EB}_k})$ . In the simulation results, we show that when  $\hat{\gamma}_{\text{CSPML}}$  differs from  $\hat{\gamma}_I$ , i.e.  $Z \nrightarrow 0$ , the impact on the variance calculation is moderately acceptable and lead to desirable results.

### 3.3 Simulation Studies

#### 3.3.1 Simulation Settings

We evaluated the performance of our proposed estimators through simulation studies in various settings and compared it to MLE from the direct regression, CSPML estimators and individual EB estimators, which incorporate single external model information. In each simulation, we compared six methods (direct regression, CSPML, EB, IVW, OCWE and SC-Learner) considering both overall and covariate-wise performance. We used the estimated standard error (ESE) to assess the precision gain of point estimates compared with the direct regression, and evaluated three overall metrics on a validation dataset of size  $N_{\text{test}}=1,000$ :

- Average estimated variance of logit-transformed predicted probability:  

$$\bar{V}[\text{logit}(\hat{p})] = \frac{1}{N_{\text{test}}} \sum_{i=1}^{N_{\text{test}}} \hat{\text{Var}}[\text{logit}(\hat{p}_i)],$$
where  $\hat{p}$  denotes the estimated probability and  $\text{logit}(p) = \log(\frac{p}{1-p}) = (\mathbf{X}, \mathbf{B})\boldsymbol{\gamma}$ ;
- Sum of squared errors:  $\text{SSE} = \frac{1}{N_{\text{test}}} \sum_{i=1}^{N_{\text{test}}} (\hat{p}_i - p_{i0})^2$ , where  $\hat{p}_{i0}$  denotes the true probability of  $Y_i = 1$  given  $X_i$  and  $B_i$ ;
- Scaled Brier score:  $\text{BS} = \frac{1}{N_{\text{test}}} \sum_{i=1}^{N_{\text{test}}} (Y_i - \hat{p}_i)^2 / \frac{1}{N_{\text{test}}} \sum_{i=1}^{N_{\text{test}}} (Y_i - \bar{Y})^2$ .

A summary of all the simulation settings is listed in Figure 3.1. In the first four simulation scenarios (I, II, III, IV), we assumed that a logistic regression model with the following form could describe the relationship among a binary outcome  $Y$  and five covariates  $(X_1, X_2, X_3, X_4, B)$ , where  $X_1, X_2, X_3, X_4$  had been used in at least one external model while  $B$  was only available in the current study:  $\text{logit}[\text{Pr}(Y = 1|X_1, X_2, X_3, X_4, B)] = -1 - 0.5 \sum_{i=1}^4 X_i + 0.5B$ . Here  $X_1, X_2, X_3, X_4$  and  $B$  followed a standard multivariate normal distribution with 0.3 correlation and the prevalence of  $Y = 1$  was 0.32.

- Simulation I evaluated the idealized case where the internal and the external models were fitted on the homogeneous population.

- Simulation II assessed the performance of the proposed method when external model 1 had biased  $\hat{\beta}$  estimates fitted from  $Y|X$ , where we obtained the incorrect model estimates by fitting external model 1 on a small dataset of size 500.
- Simulation III aims to show the impact of heterogeneous covariate distribution  $(X, B)$  in the external population. As Estes et al. [2017] assessed, the disparity of the  $(X, B)$  distribution between the two populations can come from (i) a different conditional distribution  $B|X$ , (ii) a different marginal  $B$  distribution, (iii) a different marginal  $X$  distribution, or a combination of these reasons. In this simulation scenario, we showed the combination of (i) and (ii) as an example, but we will discuss the result of other scenarios in Simulation Results in Section 3.2.
- Simulation IV evaluated the situation where the outcome model was misspecified in external model 3.

We assumed that there was an internal study of size  $n=200$  and three external models had been fitted to a very large synthetic dataset (sample size  $m_1 = m_2 = m_3 = 30,000$  for simplicity) that is sampled from the true data generating mechanism and gives precise estimates of the model parameters (except external model 1 in Simulation II). The external sample size need not be as large as 30,000 to achieve good performance as long as the estimated model parameters are close to the true parameters. Sensitivity analysis (results not shown) using external sample size  $m_1 = m_2 = m_3 = 1,000$  showed small numerical differences compared with  $m_1 = m_2 = m_3 = 30,000$ .

In simulation V and VI, we considered the outcome model with higher dimension and homogeneous populations between the internal and the external models. In these two scenarios, the internal sample size  $n=500$  was used.

- Simulation V evaluated the situation where the outcome model contained nine  $X$ 's and one  $B$ :  $\text{logit}[\Pr(Y = 1|X_1, \dots, X_9, B)] = -1 - 0.5 \sum_{i=1}^9 X_i + 0.5B$ . Specifically, external model 1 only contained two predictors,  $X_1$  and  $X_2$ , external model 2 contained seven predictors,  $X_1, \dots, X_7$ , and external model 3 contained six predictors  $X_1, X_2, X_3, X_4, X_7$  and  $X_8$ .
- simulation VI evaluated the situation where the full model contained three  $X$ 's and five  $B$ 's:  $\text{logit}[\Pr(Y = 1|X_1, X_2, X_3, B_1, \dots, B_5)] = -1 - 0.5 \sum_{i=1}^3 X_i + 0.5 \sum_{i=1}^5 B_i$ , using the same external models as simulation I.



		Internal Data	External Model 1	External Model 2	External Model 3
<b>I:</b> Correctly specified external models	Model	$Y X_1, X_2, X_3, X_4, B$	$Y X_1, X_2$	$Y X_1, X_3$	$Y X_1, X_2, X_3, X_4$
	Sample size	$n=200$	$m_1=30,000$	$m_2=30,000$	$m_3=30,000$
<b>II:</b> Biased estimate from external model 1	Model	$Y X_1, X_2, X_3, X_4, B$	$Y X_1, X_2$	$Y X_1, X_2$	$Y X_1, X_2, X_3, X_4$
	Sample size	$n=200$	$m_1=500$ (incorrect model estimates)	$m_2=30,000$	$m_3=30,000$
<b>III:</b> The joint distribution of $(X, B)$ is different in the external model 1*	Model	$Y X_1, X_2, X_3, X_4, B$	$Y X_1, X_2$	$Y X_1, X_2$	$Y X_1, X_2, X_3, X_4$
	$(X, B)$	$(X, B) \sim N((0, 0), 1)$ , correlation 0.3	$(X, B) \sim N((0, 1.5), 1)$ , correlation 0.8	Same as the internal data	Same as the internal data
	Sample size	$n=200$	$m_1=30,000$	$m_2=30,000$	$m_3=30,000$
<b>IV:</b> The outcome model is misspecified in the external model 3 †	Model	$Y X_1, X_2, X_3, X_4, B$	$Y X_1, X_2$	$Y X_1, X_3$	$Y X_1, X_2, X_3, X_4$
	Outcome model	Logit $\Pr(Y=1 X, B) = -1 - 0.5X + 0.5B$	Same as the internal data	Same as the internal data	Logit $\Pr(Y=1 X, B) = 1 + 0.5X + 0.5B$
	Sample size	$n=200$	$m_1=30,000$	$m_2=30,000$	$m_3=30,000$
<b>V:</b> homogeneous populations with correct external models (high dimensional $X$ 's)	Model	$Y X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8, X_9, B$	$Y X_1, X_2$	$Y X_1, X_2, X_3, X_4, X_5, X_6$	$Y X_1, X_2, X_3, X_4, X_7, X_8$
	Sample size	$n=500$	$m_1=30,000$	$m_2=30,000$	$m_3=30,000$
<b>VI:</b> homogeneous populations with correct external models (high dimensional $B$ 's)	Model	$Y X_1, X_2, X_3, B_1, B_2, B_3, B_4, B_5$	$Y X_1, X_2$	$Y X_1, X_3$	$Y X_1, X_2, X_3, X_4$
	Sample size	$n=500$	$m_1=30,000$	$m_2=30,000$	$m_3=30,000$

\* More scenarios that differ the joint distribution of  $(X, B)$  are presented in Appendix D of Supplementary Material.

† More scenarios that change the magnitude of outcome model are presented in Appendix D of Supplementary Material.

Note: If not specifically stated,  $(X, B)$  always followed Gaussian distribution with mean zero, standard deviation 1 and correlation 0.3; the outcome model used was the logistic model with the form  $\text{logit } \Pr(Y=1|X, B) = -1 - 0.5X + 0.5B$ .

Figure 3.1: Simulation settings snapshot

### 3.3.2 Simulation Results

In Table 3.1 simulation I, we see that all CSPML and EB estimators are unbiased as expected. The estimated standard error (ESE, in square brackets) accurately reflect the true standard deviation (SD, in round brackets) from 500 simulations. Both OCWE and SC-Learner had better overall performance than single EB estimators, while SC-Learner outperformed OCWE with respect to both the covariate-wise and the overall metrics (Figure 3.2).

Table 3.1: Results of simulation I–IV. Internal dataset had size  $n=200$ ; green represents good performance with small bias and large efficiency gain, yellow represents underestimated ESE (in square brackets) compared with SD (in round brackets), and red represents poor performance of bias/95% coverage rate.

Simulation	Direct Regression		Internal Data + External 1		Internal Data + External 2		Internal Data + External 3		Composite of EB Estimators		SC-Learner
	Bias (SD) [ESE]	95% Coverage Rate	CSPML 1	EB 1	CSPML 2	EB 2	CSPML 3	EB 3	IVW	OCWE	
I	Weight	/	/	/	/	/	/	/	(.33, .33, .33)	(.27, .12, .61)	/
	$\gamma_0$	-0.41 (.197) [.193]	-0.08 (.081) [.085]	-0.32 (.162) [.160]	-0.07 (.089) [.084]	-0.32 (.163) [.161]	.004 (.050) [.085]	-0.33 (.169) [.168]	-0.32 (.164) [.162]	-0.31 (.166) [.164]	-0.32 (.164) [.162]
		95%	96%	95%	95%	95%	98%	95%	95%	95%	95%
	$\gamma_1$	-0.20 (.229) [.208]	-0.04 (.095) [.097]	-0.14 (.190) [.173]	-0.02 (.089) [.097]	-0.13 (.188) [.174]	-0.02 (.066) [.074]	-0.16 (.198) [.181]	-0.14 (.191) [.176]	-0.16 (.194) [.178]	-0.14 (.191) [.175]
		94%	96%	93%	97%	93%	96%	93%	93%	93%	93%
	$\gamma_2$	-0.17 (.219) [.208]	-0.08 (.091) [.096]	-0.14 (.180) [.173]			-0.03 (.062) [.210]	-0.14 (.188) [.182]	-0.15 (.195) [.187]	-0.11 (.190) [.181]	-0.14 (.184) [.176]
		94%	97%	94%			97%	94%	94%	94%	94%
II	Weight	/	/	/	/	/	/	/	(.33, .33, .34)	(.23, .28, .49)	/
	$\gamma_0$	-0.42 (.198) [.193]	.042 (.080) [.086]	-0.24 (.168) [.161]	-0.08 (.081) [.085]	-0.32 (.162) [.160]	.003 (.050) [.058]	-0.34 (.170) [.169]	-0.31 (.166) [.162]	-0.31 (.167) [.163]	-0.31 (.166) [.162]
		95%	87%	94%	96%	95%	98%	95%	95%	95%	95%
	$\gamma_1$	-0.20 (.229) [.208]	.153 (.094) [.100]	-0.15 (.202) [.174]	-0.04 (.095) [.097]	-0.14 (.190) [.173]	-0.17 (.066) [.074]	-0.18 (.198) [.182]	-0.07 (.196) [.176]	-0.09 (.197) [.177]	-0.07 (.196) [.175]
		94%	65%	92%	96%	93%	97%	94%	93%	92%	93%
	$\gamma_2$	-0.18 (.220) [.208]	-143 (.091) [.098]	-0.41 (.190) [.174]	-0.08 (.091) [.096]	-0.14 (.180) [.173]	.009 (.062) [.073]	-0.13 (.188) [.183]	-0.22 (.185) [.176]	-0.20 (.186) [.177]	-0.28 (.189) [.177]
		94%	70%	93%	96%	94%	95%	94%	94%	94%	94%
III	Weight	/	/	/	/	/	/	/	(.33, .33, .33)	(.15, .35, .49)	/
	$\gamma_0$	-0.41 (.198) [.193]	.648 (.082) [.100]	-0.09 (.199) [.162]	-0.39 (.081) [.085]	-0.04 (.162) [.160]	-0.33 (.051) [.058]	-0.39 (.169) [.168]	-0.29 (.176) [.163]	-0.31 (.173) [.163]	-0.29 (.176) [.162]
		95%	0%	91%	97%	95%	98%	95%	93%	93%	93%
	$\gamma_1$	-0.19 (.228) [.208]	-156 (.096) [.113]	-0.26 (.221) [.175]	.012 (.095) [.097]	-0.10 (.189) [.173]	.022 (.066) [.074]	-0.12 (.197) [.181]	-0.16 (.202) [.176]	-0.14 (.198) [.177]	-0.16 (.202) [.176]
		94%	75%	89%	97%	93%	97%	94%	91%	93%	91%
	$\gamma_2$	-0.17 (.219) [.208]	.160 (.091) [.113]	-0.24 (.211) [.175]	-0.15 (.091) [.096]	-0.15 (.179) [.173]	-0.05 (.062) [.073]	-0.14 (.188) [.182]	-0.17 (.192) [.176]	-0.16 (.188) [.177]	-0.19 (.200) [.177]
		95%	76%	90%	97%	94%	95%	94%	93%	94%	93%
IV	Weight	/	/	/	/	/	/	/	(.34, .33, .32)	(.57, .38, .04)	/
	$\gamma_0$	-0.41 (.198) [.193]	-0.08 (.081) [.085]	-0.32 (.163) [.161]	-0.07 (.089) [.084]	-0.32 (.164) [.161]	.169 (.189) [.162]	-0.31 (.169) [.178]	-0.32 (.173) [.165]	-0.30 (.164) [.160]	-0.32 (.172) [.164]
		95%	96%	94%	95%	95%	0%	96%	93%	93%	93%
	$\gamma_1$	-0.20 (.230) [.208]	-0.05 (.095) [.097]	-0.14 (.190) [.174]	-0.02 (.089) [.097]	-0.14 (.189) [.174]	.747 (.178) [.177]	-0.16 (.197) [.192]	-0.15 (.201) [.178]	-0.13 (.191) [.173]	-0.15 (.200) [.178]
		94%	96%	93%	97%	93%	3%	95%	92%	93%	92%
	$\gamma_2$	-0.18 (.220) [.208]	-0.08 (.091) [.096]	-0.14 (.180) [.175]			.756 (.174) [.177]	-0.14 (.188) [.192]	-0.15 (.206) [.189]	-0.13 (.196) [.187]	-0.14 (.197) [.177]
		95%	97%	94%			3%	94%	93%	94%	93%
V	Weight	/	/	/	/	/	/	/	(.34, .33, .32)	(.57, .38, .04)	/
	$\gamma_0$	-0.41 (.198) [.193]	-0.08 (.081) [.085]	-0.32 (.163) [.161]	-0.07 (.089) [.084]	-0.32 (.164) [.161]	.762 (.167) [.175]	-0.14 (.178) [.191]	-0.15 (.194) [.187]	-0.13 (.193) [.189]	-0.14 (.187) [.181]
		95%	96%	94%	95%	95%	1%	95%	95%	94%	95%
	$\gamma_1$	-0.20 (.230) [.208]	-0.05 (.095) [.097]	-0.14 (.190) [.174]	-0.02 (.089) [.097]	-0.14 (.189) [.174]	.746 (.179) [.176]	-0.16 (.197) [.192]	-0.25 (.216) [.202]	-0.26 (.217) [.207]	-0.26 (.216) [.193]
		94%	96%	93%	97%	93%	4%	93%	94%	95%	94%
	$\gamma_2$	-0.18 (.220) [.208]	-0.08 (.091) [.096]	-0.14 (.180) [.175]				-0.22 (.188) [.193]	-0.25 (.216) [.202]	-0.26 (.217) [.207]	-0.26 (.216) [.193]
		95%	97%	94%				93%	94%	95%	94%

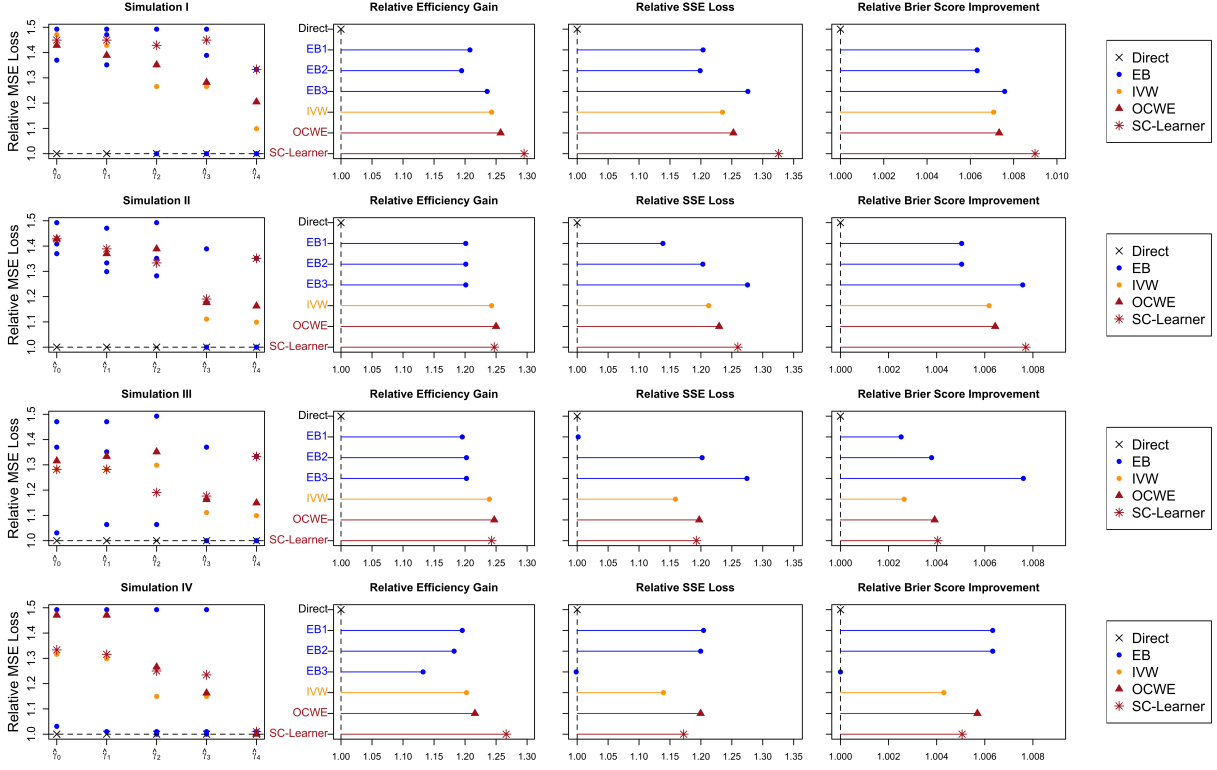


Figure 3.2: Visualization of the metrics to evaluate the performance in simulations I-IV. Scatter plot shows the covariate-wise relative MSE improvement compared with the direct regression (x axis represents  $\hat{\gamma}_X$ ;  $\hat{\gamma}_B$  not shown since no external information incorporated; larger y values represent larger MSE loss) while the line plots represent the relative efficiency/SSE/BS improvement compared with the direct regression fitting on a validation dataset of size 1,000 (longer lines represent larger improvement).

In Table 3.1 simulation II,  $\text{CSPML}_1$  was biased with poor 95% coverage rate due to biased estimation of  $\hat{\beta}$ . Although trading off most of the precision gain,  $\text{EB}_1$  corrected the bias in  $\text{CSPML}_1$ . The bias of  $\text{CSPML}_1$  also caused underestimation of the standard error of  $\text{EB}_1$  as highlighted in yellow, which was because  $\hat{\gamma}_{\text{CSPML}_1}$  differed from  $\hat{\gamma}_1$  as pointed out in Section 3.2.4. Despite that, both OCWE and SC-Learner detected the least efficiency gain of  $\text{EB}_3$  among all EB estimators and provided unbiased overall estimates. Moreover, OCWE and SC-Learner had the largest relative efficiency gain compared with the direct regression than other methods, as well as the best performance regarding the relative SSE and BS gain only second to  $\text{EB}_3$  (Figure 3.2).

Results of simulation III and IV in Table 3.1 indicates that when incorporating discordant external information from heterogeneous populations (one with different distribution of the joint  $[X, B]$  and another with different outcome model), EB estimators are able to correct the bias of CSPML estimators by trading off some precision gain, and the proposed estimators could further identify and down-weight that particular EB estimators. Estes et al. [2017] provided comprehensive sim-

ulation results to show that the EB estimator can protect against the bias due to heterogeneous  $(X, B)$  distribution from the external model. Their results indicated that CSPML estimator would have a substantial bias when the difference came from the conditional distribution  $B|X$  or marginal distribution  $B$ , and marginal distribution  $X$  when the  $X$ - $B$  interaction term is involved in the true  $Y|X, B$  model. When evaluating these scenarios (summary and results in Appendix B.1), we found a similar pattern and conclusion as simulation III. In summary, these simulation results provided evidence that the proposed approach is robust to the heterogeneous covariate distribution  $(X, B)$  of the external population.

In addition to simulation IV, we have assessed different forms of misspecified outcome model in the external population listed in Appendix B.4, including different intercept only and different  $X$  coefficients only, which showed similar conclusions. Simulation IV also showed that although the higher dimension of the external model will lead to a better overall prediction when the internal and external has the same population (simulation I and II), it is not the case when the transportability assumption is violated: OCWE identified that external model 3 came from a different distribution than the internal study and thus assigned the smallest weight to  $EB_3$ , even though it had the greatest number of predictors compared with the other two external models. In both simulations, OCWE and SC-Learner had similar covariate-wise and overall performances (Figure 3.2).

Table 3.2 further shows that the proposed estimators have decent performance when the number of predictors in the external models has large differences (simulation V) and when the dimension of  $B$  is larger than  $X$ 's (simulation VI). In simulation V, compared with other external models that included more than six predictors, external model 1 only used two predictors and thus provided the least amount of extra information. We see that OCWE assigns the minimum weight to  $EB_1$  and achieves the largest relative efficiency gain among all estimators (Figure 3.3). The fact that SC-Learner outperformed OCWE with respect to decreasing the covariate-wise variance of  $\hat{\gamma}_5$  to  $\hat{\gamma}_9$  reveals that SC-Learner is flexible enough to select the external information on the covariate-level when  $p_k$  (the number of predictors used in the  $k^{\text{th}}$  external model) is very different from others. Simulation VI showed that when the dimension of  $B$  was much larger than  $X$ , the overall benefit of combining multiple EB estimators would be limited due to the small amount of external information added from the external models.

Table 3.2: Results of simulation V–VI. Internal dataset had size  $n=500$ ; green represents good performance with small bias and large efficiency gain.

Simulation	Direct Regression		Internal Data + External 1			Internal Data + External 2			Internal Data + External 3			Composite of EB Estimators			SC-Learner		
	Bias (SD) [ESE]	95% Coverage Rate	CSPML 1			CSPML 2			CSPML 3			EB 3				IVW	OCWE
V	Weight	/	/	/	/	/	/	/	/	/	/	/	(.33, .33, .33)	(.13, .432, .429)	/		
	$\gamma_0$	-0.35 (.153) [.149] 94%	.02 (.100) [.098] 95%	.031 (.132) [.129] 95%	-0.02 (.054) [.059] 97%	-0.031 (.139) [.137] 94%	-0.10 (.064) [.066] 96%	-0.031 (.137) [.135] 94%	-0.031 (.135) [.132] 94%	-0.031 (.137) [.133] 94%	-0.031 (.137) [.133] 94%	-0.031 (.137) [.133] 94%	-0.031 (.135) [.132] 94%	-0.031 (.135) [.132] 94%	-0.031 (.135) [.132] 94%		
	$\gamma_1$	-0.23 (.160) [.158] 94%	.002 (.089) [.094] 96%	.016 (.133) [.133] 94%	.010 (.055) [.061] 97%	-0.019 (.144) [.144] 94%	.009 (.061) [.067] 97%	-0.019 (.142) [.142] 94%	-0.018 (.139) [.139] 94%	-0.019 (.141) [.141] 94%	-0.019 (.141) [.141] 94%	-0.019 (.141) [.141] 94%	-0.018 (.139) [.139] 94%	-0.018 (.139) [.139] 94%	-0.018 (.139) [.139] 94%		
	$\gamma_2$	-0.27 (.153) [.158] 95%	.014 (.094) [.094] 95%	.023 (.128) [.133] 95%	-0.006 (.056) [.060] 97%	-0.024 (.137) [.144] 94%	.013 (.064) [.066] 96%	-0.025 (.136) [.142] 95%	-0.024 (.133) [.139] 95%	-0.024 (.135) [.141] 95%	-0.024 (.135) [.141] 95%	-0.024 (.135) [.141] 95%	-0.024 (.133) [.139] 95%	-0.024 (.133) [.138] 95%	-0.024 (.133) [.138] 95%		
	$\gamma_3$	-0.37 (.161) [.158] 95%			.002 (.057) [.061] 97%	-0.031 (.144) [.145] 95%	.002 (.064) [.067] 97%	-0.031 (.143) [.143] 95%	-0.033 (.149) [.148] 95%	-0.031 (.145) [.145] 95%	-0.031 (.145) [.145] 95%	-0.031 (.145) [.145] 95%	-0.031 (.143) [.143] 95%	-0.031 (.143) [.143] 95%	-0.031 (.143) [.143] 95%		
	$\gamma_4$	-0.14 (.162) [.157] 94%			.025 (.054) [.060] 97%	-0.014 (.145) [.144] 95%	.025 (.062) [.066] 96%	-0.016 (.144) [.142] 95%	-0.015 (.150) [.147] 95%	-0.014 (.147) [.144] 95%	-0.014 (.147) [.144] 95%	-0.014 (.147) [.144] 95%	-0.015 (.145) [.142] 96%	-0.015 (.145) [.142] 96%	-0.015 (.145) [.142] 96%		
	$\gamma_5$	-0.14 (.160) [.158] 96%			.010 (.055) [.060] 97%	-0.011 (.143) [.153] 95%				-0.013 (.154) [.150] 95%	-0.012 (.153) [.150] 95%	-0.012 (.153) [.150] 95%	-0.011 (.143) [.153] 96%	-0.011 (.143) [.153] 96%	-0.011 (.143) [.153] 96%		
	$\gamma_6$	-0.23 (.160) [.158] 95%			-0.010 (.057) [.061] 99%	-0.021 (.143) [.153] 96%				-0.023 (.154) [.151] 95%	-0.022 (.152) [.151] 95%	-0.022 (.152) [.151] 95%	-0.021 (.143) [.153] 96%	-0.021 (.143) [.153] 96%	-0.021 (.143) [.153] 96%		
	$\gamma_7$	-0.15 (.165) [.158] 95%			-0.008 (.056) [.060] 99%	-0.014 (.148) [.153] 96%				-0.015 (.159) [.150] 95%	-0.015 (.157) [.150] 95%	-0.015 (.157) [.150] 95%	-0.014 (.148) [.153] 96%	-0.014 (.148) [.153] 96%	-0.014 (.148) [.153] 96%		
	$\gamma_8$	-0.33 (.159) [.159] 95%								-0.001 (.063) [.067] 98%	-0.029 (.141) [.151] 96%	-0.031 (.153) [.152] 97%	-0.031 (.151) [.151] 97%	-0.031 (.151) [.151] 97%	-0.029 (.141) [.151] 97%	-0.029 (.141) [.151] 97%	
$\gamma_9$	-0.27 (.154) [.154] 97%								-0.021 (.066) [.067] 97%	-0.026 (.136) [.151] 97%	-0.026 (.148) [.151] 97%	-0.026 (.146) [.150] 97%	-0.026 (.146) [.150] 97%	-0.026 (.136) [.151] 97%	-0.026 (.136) [.151] 97%		
$\gamma_B$	.034 (.172) [.172] 97%																
VI	Weight	/	/	/	/	/	/	/	/	/	/	/	[.33, .33, .34]	[.01, .0004, .987]	/		
	$\gamma_0$	-0.27 (.126) [.121] 94%	.003 (.078) [.073] 92%	-0.020 (.108) [.103] 94%	.003 (.078) [.073] 93%	-0.031 (.108) [.103] 94%	.005 (.078) [.073] 92%	-0.020 (.110) [.105] 94%	-0.020 (.109) [.104] 94%	-0.020 (.110) [.105] 94%	-0.020 (.110) [.105] 94%	-0.020 (.110) [.105] 94%	-0.020 (.109) [.104] 94%	-0.020 (.109) [.104] 94%	-0.020 (.109) [.104] 94%		
	$\gamma_1$	-0.12 (.132) [.132] 95%	-0.015 (.088) [.085] 95%	-0.013 (.115) [.115] 95%	-0.019 (.086) [.085] 94%	-0.019 (.115) [.115] 95%	-0.019 (.084) [.081] 95%	-0.013 (.116) [.116] 95%	-0.013 (.115) [.115] 95%	-0.013 (.116) [.116] 95%	-0.013 (.116) [.116] 95%	-0.013 (.115) [.115] 95%	-0.013 (.115) [.115] 95%	-0.013 (.115) [.115] 95%	-0.013 (.115) [.115] 95%		
	$\gamma_2$	-0.04 (.137) [.132] 93%	.005 (.087) [.085] 95%	.004 (.118) [.115] 93%			.004 (.084) [.081] 94%	.003 (.119) [.117] 94%	.004 (.124) [.121] 93%	.003 (.119) [.117] 94%	.003 (.119) [.117] 94%	.004 (.118) [.116] 94%	.004 (.118) [.116] 94%	.004 (.118) [.116] 94%	.004 (.118) [.116] 94%		
	$\gamma_3$	-0.11 (.131) [.132] 96%			.005 (.085) [.085] 95%	-0.007 (.113) [.117] 96%	.005 (.081) [.082] 94%	-0.007 (.114) [.118] 96%	-0.008 (.119) [.121] 96%	-0.007 (.114) [.118] 96%	-0.007 (.114) [.118] 96%	-0.007 (.113) [.117] 96%	-0.007 (.113) [.117] 96%	-0.007 (.113) [.117] 96%	-0.007 (.113) [.117] 96%		
	$\gamma_{B1}$	.007 (.126) [.121] 95%															
	$\gamma_{B2}$	-0.25 (.137) [.132] 96%															
	$\gamma_{B3}$	-0.14 (.135) [.132] 97%															
	$\gamma_{B4}$	-0.12 (.131) [.132] 96%															
	$\gamma_{B5}$	-0.13 (.139) [.132] 95%															

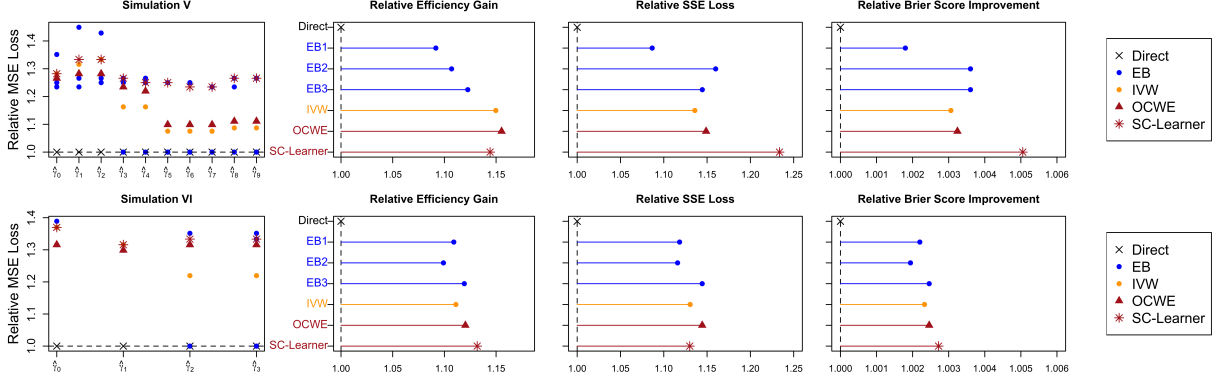


Figure 3.3: Visualization of the metrics to evaluate the performance in simulations V-VI. Scatter plot shows the covariate-wise relative MSE improvement compared with the direct regression (x axis represents  $\hat{\gamma}_X$ ;  $\hat{\gamma}_B$  not shown since no external information incorporated; larger y values represent larger MSE loss) while the line plots represent the relative efficiency/SSE/BS improvement compared with the direct regression fitting on a validation dataset of size 1,000 (longer lines represent larger improvement).

### 3.4 Application to Prostate Cancer Data

To assess the performance of the proposed estimators in a real data example, we developed a model for predicting the risk of high-grade prostate cancer (Gleason score  $> 6$ ) using a combination of a set of internal individual-level data and two external risk calculators from different studies. The first risk calculator was developed based on the Prostate Cancer Prevention Trial in the United States [Thompson et al., 2006]. This calculator, denoted as PCPThg, is built on five clinical variables including prostate-specific antigen (PSA) level, digital rectal examination (DRE) findings, age, race (African American or not) and prior biopsy results using the following model:

$$\text{logit}(p_i) = -6.25 + 1.29\log(\text{PSA}_i) + \text{DRE}_i + 0.03\text{Age}_i + 0.96\text{Race}_i - 0.36\text{Biopsy}_i, \quad (3.3)$$

where  $p_i$  is the probability of observing high-grade prostate cancer for subject  $i$ . The second risk calculator is the European Randomized Study of Screening for Prostate Cancer (ERSPC) risk calculator 3 [Roobol et al., 2012], which uses slightly different clinical variables to predict the same risk as PCPThg: PSA level, DRE findings, and transrectal ultrasound prostate volume (TRUS-PV) in a logistic regression model shown as below:

$$\text{logit}(p_i) = -3.15 + 1.18\log_2(\text{PSA}_i) + 1.81\text{DRE}_i - 1.51\log_2(\text{TRUS-PV}_i), \quad (3.4)$$

where TRUS-PV was categorized as a 3-category variable described in Roobol et al. [2012]. In addition to all the predictors used in the external model 3.3 and 3.4, we considered adding two more log-2-transformed biomarkers that had not been widely used but shown to be predictive of prostate cancer [Tomlins et al., 2015, Truong et al., 2013], prostate cancer antigen 3 (PCA3) and TMPRSS2:ERG (T2:ERG) gene fusions to our target model:

$$\begin{aligned} \text{logit}(p_i) = & \gamma_0 + \gamma_1 \log_2(\text{PSA}_i) + \gamma_2 \text{DRE}_i + \gamma_3 \text{Age}_i + \gamma_4 \text{Biopsy}_i + \gamma_5 \text{Race}_i \\ & + \gamma_6 \log_2(\text{TRUS-PV}_i) + \gamma_7 \log_2(\text{PCA3}_i + 1) + \gamma_8 \log_2(\text{T2:ERG}_i + 1), \end{aligned} \quad (3.5)$$

Using the data from Tomlins et al. [2015], we had 678 male patients in the internal dataset who had complete data of all eight covariates listed in model 3.5 and an independent validation data of size 1,174 (sample size reduced from the initial 679 and 1218 patients, respectively due to missing TRUS-PV that was not previously used in Tomlins et al. [2015]). Details of the individual-level data, including the description of the internal and the validation dataset, and the recent applications using the same setting can be found in Tomlins et al. [2015] and Cheng et al. [2019].

PCPThg utilizes standard clinical and demographic variables that have been widely used while ERSPC additionally considers the prostate volume that was shown to be related to PSA level [Bohnen et al., 2007]. In addition, similarities of prostate-specific antigen patterns between the United States and European populations prostate-specific antigen patterns have been shown [Simpkin et al., 2016]. Therefore, incorporating information from both risk calculators can potentially provide more accurate estimation of the risk parameters and narrower confidence bands, which could in turn yield better prediction performance and improved inference.

In order to reconcile the discrepancy between the covariates used in the external models and have a compatible interpretation of the intercept, we centered all variables in the original models 3.3 and 3.4, and log2-transformed PSA and TRUS-PV by adjusting the corresponding intercepts a-priori (details in Appendix B.5). We present the estimated coefficients and standard errors in Table 3.3. Similar to the simulation study, we calculated the scaled Brier Score and the average estimated variance of logit-transformed predicted probability based on the validation dataset as the prediction metrics.

As shown in Table 3.3, OCWE assigns almost zero weight to the ERSPC, which indicated a large population discrepancy between the internal data and the underlying European population possibly due to the difference in the intercept, which reflects that the prevalence of high-grade prostate cancer is higher in the European population compared with patients in the United States who had average covariate values. Even though, SC-Learner was able to make the most of the little improvement provided by ERSPC and augment the point-wise precision gain for covariate PSA, DRE and TRUS-PV (3%, 4% and 12% more compared with OCWE, respectively), which



led to the largest overall improvement as well (17.2 % decrease of the average prediction variance compared with the direct regression).

### 3.5 Discussion

The proposed framework along with two weighted estimators, OCWE and SC-Learner, adds to the evolving research on using external summary-level information to bolster the statistical efficiency of the internal study for improved inference. This new method is flexible and robust in the ways that (i) it is capable of incorporating external models that use a slightly different set of covariates; (ii) it is able to identify the most relevant external information and diminish the influence of information that is less compatible with the internal data; and (iii) it nicely balances the bias-variance trade off while preserving the most precision gain. Moreover, our extensive simulation studies and the real data example show that the proposed estimators are more efficient and robust than the naïve analysis of the internal data and other naïve combinations of external estimators in both idealized and non-idealized settings.

Compared with a single EB estimator, the proposed composite estimators can have up to 32% more improvement in MSE regarding one covariate (Figure 3.2) and decent improvement regarding the overall metric such as 20% further decrease in SSE and 11.5% further decrease in estimated prediction variance (Figure 3.2). In some cases, several single EB estimators that showed limited gains mitigate another single EB estimator’s excellent performance during integration, e.g. Simulation II. In reality, the proposed composite estimators will be preferred over a single EB estimator, since it is often difficult to pick the best external model that contains the most useful information to boost the inference of the internal study among several available external models.

In practice, the choice of SC-Learner or OCWE mainly depends on the features of the external models and the user’s research goal. As shown in simulation V, if at least one of the external models used very few predictors compared to the full dimension of  $(X, B)$ , i.e.,  $p_k \ll p$ , we suggest using SC-Learner as it can adapt the external information being considered covariate-wise and thus prevent the gain in certain covariates from being washed away when the dimension of predictors are uneven. Similarly, we would recommend using SC-Learner if the researcher cares about maximizing the precision gain on the covariate-level or improving precision in certain covariates are of particular interest. On the other hand, if the research goal involves ranking the usefulness/relevancy of the external model, OCWE would be a good choice as it provides one unified weight for all covariates in the same model and outperforms comparable estimators such as IVW (inverse variance-weighted estimator).

A different approach to estimating the weights in OCWE would be cross-validation. For example, instead of minimizing the estimated prediction variance over all internal observations, one



Table 3.3: Results of the real data example for predicting the risk of high-grade prostate cancer. Internal dataset of size  $n=678$ ; validation dataset of size  $N_{\text{test}}=1,174$ ;  $\downarrow$  %, percentage of ESE decrease compared with the direct regression; Firth correction applied in direct regression, and estimated PCPThg and ERSPC; REF, reference; green represents good performance with small bias and large efficiency gain, and grey represents no efficiency improvement due to no added information from the external calculator.

	Direct Regression			PCPThg			ERSPC			Internal data + PCPThg			Internal data + ERSPC			Composite of EB Estimators		
	Original	Estimated	Original	Estimated	Original	Estimated	Original	Estimated	Original	CSPML 1	EB 1	PCPThg	CSPML 2	EB 2	ERSPC	IVW	OCWE	SC-Leamer
Weight	/	/	/	/	/	/	/	/	/	/	/	/	/	/	/	[.5, .5]	[.99, .00]	/
$\bar{V}[\logit(\hat{p})]$	.1199	/	/	.065	/	.048	/	.073	.078	.104	.105	.104 ( $\downarrow$ 13.2%)	.104 ( $\downarrow$ 13.2%)	.105	.104 ( $\downarrow$ 13.2%)	.104 ( $\downarrow$ 13.2%)	.104 ( $\downarrow$ 13.2%)	.0996 ( $\downarrow$ 17.2%)
Brier Score	.902	1.059	.987	.994	.954	.994	.942	.942	.901	.901	.901	.901	.901 ( $\downarrow$ .12%)	.901	.901 ( $\downarrow$ .12%)	.901 ( $\downarrow$ .12%)	.901 ( $\downarrow$ .12%)	.901 ( $\downarrow$ .12%)
Point estimate (ESE)																		
$\downarrow$ % ESE w.r.t direct regression																		
Intercept	-4.123 (.449)	-3.686268	-1.395 (.115)	-3.16	-1.382 (.111)	-6.422 (.437)	-4.128 (.445)	-5.989 (.437)	-4.130 (.444)	-4.129 (.443)	-4.128 (.445)	-4.129 (.443)	-4.128 (.445)	-4.129 (.443)	-4.128 (.445)	-4.129 (.443)	-4.128 (.445)	-4.129 (.443)
REF	REF	/	74%	/	75%	3%	1%	3%	1%	1%	1%	1%	1%	1%	1%	1%	1%	1%
$\log_2(\text{PSA})$	0.860 (.144)	0.8941599	0.721 (.123)	1.175573	.878 (.123)	0.893 (.103)	0.860 (.132)	1.159 (.081)	0.861 (.124)	0.860 (.127)	0.860 (.130)	0.860 (.127)	0.860 (.130)	0.861 (.124)	0.860 (.130)	0.860 (.127)	0.860 (.130)	0.860 (.127)
REF	REF	/	14%	/	14%	28%	8%	44%	14%	11%	9%	11%	9%	14%	11%	12%	9%	12%
DRE	1.028 (.297)	1	1.134 (.256)	1.813195	1.298 (.269)	0.687 (.218)	1.027 (.269)	1.430 (.170)	1.029 (.253)	1.028 (.259)	1.027 (.269)	1.028 (.259)	1.027 (.269)	1.029 (.253)	1.028 (.259)	1.028 (.258)	1.027 (.269)	1.028 (.258)
REF	REF	/	14%	/	9%	26%	9%	43%	15%	13%	9%	13%	9%	15%	13%	13%	9%	13%
Age	0.0315 (.014)	0.03	0.0328 (.012)			0.0304 (.009)	0.0315 (.013)	0.032 (.014)	0.031 (.014)	0.032 (.014)	0.032 (.013)	0.032 (.013)	0.032 (.013)	0.031 (.014)	0.032 (.013)	0.032 (.013)	0.032 (.013)	0.032 (.013)
REF	REF	/	14%	/		31%	8%	31%	1%	1%	8%	1%	8%	1%	8%	8%	8%	8%
Biopsy	-1.165 (.290)	-0.36	-1.413 (.270)			-0.0286 (.158)	-1.163 (.255)	-1.187 (.290)	-1.165 (.291)	-1.164 (.272)	-1.163 (.255)	-1.163 (.255)	-1.163 (.255)	-1.165 (.291)	-1.163 (.255)	-1.163 (.255)	-1.163 (.255)	-1.163 (.255)
REF	REF	/	7%	/		46%	12%	46%	6%	6%	12%	6%	12%	6%	12%	12%	12%	12%
Race	0.193 (.329)	0.96	0.448 (.287)			0.819 (.218)	0.194 (.291)	0.160 (.329)	0.193 (.329)	0.193 (.308)	0.194 (.291)	0.194 (.291)	0.194 (.291)	0.193 (.329)	0.194 (.291)	0.194 (.291)	0.194 (.291)	0.194 (.291)
REF	REF	/	13%	/		34%	12%	34%	6%	6%	12%	6%	12%	6%	12%	12%	12%	12%
$\log_2(\text{TRUS-PV})$	-1.663 (.252)					-1.683 (.252)	-1.663 (.252)	-1.491 (.155)	-1.663 (.223)	-1.663 (.236)	-1.663 (.252)	-1.663 (.236)	-1.663 (.252)	-1.663 (.223)	-1.663 (.252)	-1.663 (.223)	-1.663 (.252)	-1.663 (.223)
REF	REF	/				11%		38%	12%	6%	12%	6%	0%	12%	0%	12%	0%	12%
$\log_2(\text{PCA3+1})$	.485 (.088)					0.495 (.088)	0.485 (.088)	0.507 (.088)	0.485 (.088)	0.485 (.088)	0.485 (.088)	0.485 (.088)	0.485 (.088)	0.485 (.088)	0.485 (.088)	0.485 (.088)	0.485 (.088)	0.486 (.088)
REF	REF	/																
$\log_2(\text{T2-ERG+1})$	.096 (.037)					0.095 (.037)	0.096 (.037)	0.099 (.037)	0.096 (.037)	0.096 (.037)	0.096 (.037)	0.096 (.037)	0.096 (.037)	0.096 (.037)	0.096 (.037)	0.096 (.037)	0.096 (.037)	0.096 (.037)
REF	REF	/																

could randomly split the data into training and testing data, and choose the weights that minimize the objective function over the testing data only, then average this over different data splits. This approach could potentially prevent overfitting, give more stable estimates of the weights and improve the predictive performance.

As is typical of shrinkage estimators, in finite-sized samples, the EB method sacrifices a small amount of efficiency when the assumption of full transportability is satisfied, but reduces the potential bias of the CSPML when full transportability is not satisfied, while still being more efficient than the simple estimate from the internal data. In the proposed method, the magnitude of the precision gain depends on the degree of the distributional similarity between the internal and external populations, i.e., the more similar, the more benefit we will gain by incorporating the external models. In the extreme case when these populations are completely different, our approach is very similar to analyzing the internal data only. On the contrary, when these populations share the identical joint distribution of  $(Y, X, B)$ , we will achieve near to the maximum possible benefit.

Note that the proposed framework is not suitable if some predictors used in the external models are completely unmeasured in the internal study. In addition, the proposed framework is constructed based on parametric regression models, which requires the exact form of the external models and common covariates shared across different external models. In some cases, such as the real data example in this study, the authors were able to reconcile the discrepant transformations (i.e. one used natural-log PSA and the other used mean-centered log-2 based PSA) by reparametrizing the intercept. But this may not always be feasible, in which case we suggest considering methods that only require the predicted probability or the ability to estimate the probability given predictors without knowing the exact models, such as the synthetic data method proposed in Chapter 2. It is plausible that the assumed parametric model is not a good approximation of the internal data's underlying distribution. In the simulation (results not shown here), we saw that when the effect of non-linear terms was small, the proposed method could still correctly estimate the main effect. Besides, some reassurance about the selected structure of the parametric internal model can be obtained from the external models. The external models determine the  $X$  variables that are included and how they are included. For example, in our prostate cancer example, the external models took a log transformation of PSA, thus the parametric model for the internal data also includes a log transformation of PSA. Since the external datasets are typically large, we might surmise that if a large non-linearity or a strong interaction amongst the  $X$  variables were needed, it would have been included in the external model.

Recently, Zhang et al. [2020] proposed a general framework as the extension of Chatterjee et al. [2016] to solve the same genre of problem when the external parameter uncertainty cannot be ignored. In the situation where the external study population differs from the internal one, the performance of their method depends on the availability of high-quality reference data from the

external population, similar to in Chatterjee et al. [2016]. Our method also provides the option of incorporating external parameter uncertainty but more importantly it provides valid internal inference when it is in general hard to obtain the right reference data in reality. In addition, Kundu et al. [2019] proposed a generalized meta-analysis framework building from the Chatterjee et al. [2016] approach to combine information of multiple regression models with disparate covariates using a method of moment approach. Different from our goal, their method is an extension of the fixed-effect meta-analysis that also relies on the existence of reference data, and the performance of the proposed estimator depends on the quality and representativeness of the reference data.

One possible extension of the proposed method is the application in causal inference. If the treatment indicator was available as one of the  $X$  covariates, one could directly calculate the estimated average causal effect through the formula  $E_{X,B}[E(Y|\text{treatment} = 1, X, B) - E(Y|\text{treatment} = 0, X, B)]$  using the regression estimates obtained from the proposed method. Yang and Ding [2020a] considered a similar setting, where they view the internal data as the validation data with richer covariates while the external data serves as the main dataset with fewer covariates, aiming to improve the efficiency of the initial estimator  $\hat{\gamma}_I$  from the internal dataset by incorporating a constructed zero-mean error-prone estimator  $\hat{\beta}_I - \hat{\beta}_E$ , where  $\hat{\beta}_I$  and  $\hat{\beta}_E$  are the estimators using  $X$  only from the internal and external population, respectively. Using notation from Estes et al. [2017], we can also reparametrize Yang and Ding’s estimator as a weighted average of  $\hat{\gamma}_I$  and  $\hat{\gamma}_{\text{CSPML}}$ , where the only difference is that the EB estimator has the shrinkage effect by empirically estimating the variance-covariance matrix that plays an important role in the weights.

It is worth noting that there is a popular field in machine learning called ensemble learning with a large and evolving literature, aiming to combine several base models to produce the optimal predictive model. Some representative ensemble methods include but not limited to Boosting [Schapire, 1990], Bagging [Freund and Schapire, 1997] and Stacking [Breiman, 1996] with some examples being random forest [Breiman, 2001] and Super Learner [van der Laan et al., 2007]. The key difference of our proposed method is that we have a specific parametric model of interest, and we are taking the weighted average of the estimated coefficients of that model from several estimators such that we can measure the impact of each predictor and its uncertainty, instead of directly weighting the predicted outcomes as in these ensemble methods. The proposed method can provide competitive and robust estimators for statistical inference. Moreover, the proposed estimators have improved efficiency compared with direct regression using the internal data only or naïve inverse variance-weighted estimator. However, if the research goal is to find the optimal predictive model with the minimum prediction error solely, especially when the underlying mechanism is not of interest, it would be worthwhile to explore the field of ensemble methods mentioned above, which is beyond the discussion of this study.

Last but not least, the issue of transportability of risk prediction models is a critical one and

one that is often encountered in practice. While it is plausible that the association between pairs of variables or even the joint distribution of all the variables is similar between populations, it is also plausible that they could differ, not just due to biological or behavioral differences in populations but also due to being collected in different parts of the world or different decades. The EB strategy will be a good choice, balancing between bias and efficiency when one is unsure about whether transportability assumptions hold for risk models across time, space or cohorts.

## 3.6 Software and Publication

R package MetaIntegration is available at <https://github.com/umich-biostatistics/MetaIntegration>. The content of this chapter has been published in *Biostatistics* at DOI: 10.1093/biostatistics/kxab017.

## CHAPTER 4

# Regression Inference for Multiple Populations by Integrating Summary-Level Data using Stacked Imputations

### 4.1 Introduction

In Chapter 3, we proposed a meta-inference framework using an empirical Bayes approach, which adaptively weighted the estimates from multiple external models to incorporate the most compatible information with the internal data while balanced the bias-variance trade-off. In this chapter, we revisit and extend the synthetic data method proposed in Chapter 2, and further allow for heterogeneity of covariate effects across the external populations, one of the key challenges in data integration and ignoring which would lead to potential estimation bias and misleading inference.

Efforts have been made to address the issue of population difference across studies. The meta-inference framework proposed in Chapter 3 assigns larger weights to the more compatible external data sources to incorporate valid supplementary information into the internal study. Chen et al. [2020] used a penalty function to identify the difference of aggregate information among data sources. Yang and Ding [2020b] employed a sensitivity parameter to quantify such systematic differences. However, the facts that different external studies may use different subsets of covariates and the underlying prediction model may be parametric or constructed by machine learning approaches add to the difficulties when making valid inference of both the internal and the external populations in data integration, let alone the external summary-level information may contain estimated regression coefficients or fitted predictions.

Although some of the existing approaches have considered heterogeneity across data sources, the main focus has been on improving the statistical efficiency of the internal dataset with little attempts to make statistical inference on external populations or allowing heterogeneous covariate effects across data sources. Wang, Wang, and Song [2012a] proposed a joint estimating procedure to merge longitudinal datasets while allowing different study-specific coefficients. The meta-

analysis approach proposed by Kundu, Tang, and Chatterjee [2019] used a generalized method of moments to estimate study-specific effects but only allows covariates that were measured in at least one of the external studies. Antonelli, Zigler, and Dominici [2017] proposed a unified Bayesian imputation framework built upon the work by Wang, Parmigiani, and Dominici [2012b], taking into account the prior odds of including a predictor in the outcome model given that it is in the exposure model, and allowing heterogeneous treatment effects by positing different population indicators in the outcome model.

In this chapter, we consider the situation where moderately sized individual data is available from the internal study, and there are  $K$  populations ( $K \geq 1$ ), each of which provides some information about the relationship between the same outcome and a slightly different set of predictors. We propose an imputation-based methodology where the goal is to fit an outcome regression model with all available variables in the internal study while utilizing summary information from external models that may have used only a subset of the predictors. The method allows for heterogeneity of covariate effects across the external populations, by first generating synthetic outcome data in each population, then using stacked multiple imputation to create a long dataset with complete covariate information, and finally analyzing the imputed data with weighted regression. This flexible and unified approach attains the following four objectives: (i) incorporating supplementary information from a broad class of externally fitted predictive models or established risk calculators based on parametric regression or machine learning methods, as long as the external model can generate outcome values given covariates; (ii) improving statistical efficiency of the estimated coefficients in the internal study; (iii) improving predictions by utilizing even partial information available from prediction models that uses a subset of the full set of covariates used in the internal study; and (iv) providing valid statistical inference for the external population with potentially different covariate effects from the internal population.

The rest of this chapter is organized as follows: In Section 4.2, we introduce the proposed methodology. In Section 4.3, we evaluate the performance of our proposed approach in a simulation study. In Section 4.4, we apply the proposed strategy to a data example, where we build an expanded risk model to predict high-grade prostate cancer borrowing information from two existing risk prediction models. Concluding remarks are presented in Section 4.5.

## 4.2 Models and Methods

### 4.2.1 Notation

Let  $Y$  denote the outcome variable of interest, either continuous or binary. Consider  $\mathbf{X}$  a set of  $P$  routinely measured variables and  $\mathbf{B}$  a set of  $Q$  new variables, e.g. newly discovered biomarkers,

where  $\mathbf{B}$  is only available in the internal study (i.e.  $\mathbf{B}$  are unmeasured variables in all of the external studies). Let  $S_k$ ,  $k=0, \dots, K$ , denote the indicators of the  $K+1$  study populations, where  $S_0$  represents the internal population and  $S_k$ 's represents the  $K$  external populations.

We assume that a moderate dataset of size  $n$  with complete variables  $Y$ ,  $\mathbf{X}$  and  $\mathbf{B}$  is available to us from the internal study. For each external population  $k \geq 1$ , a well-established reduced model for the same outcome  $Y$  is also available, each of which may use a slightly different set of predictors  $\mathbf{X}_k$ , a subset of  $\mathbf{X}$ . For example, if linear regression is used, the prediction model may look like:  $E(Y|\mathbf{X}_k) = \beta_0 + \mathbf{X}_k\beta_k$ , where the dimension of  $\mathbf{X}_k$  is  $P_k \leq P$ . We do not have access to the underlying individual-level data that was used to fit the external model but only the summary information. This summary information can come in different forms that we summarize into two categories:

**Category 1:** Directly available in the form of an externally fitted parametric regression model, along with the estimated model parameters  $\hat{\beta}_k$ ;

**Category 2:** Any parametric or non-parametric models without knowing the exact form e.g. established risk calculators that provide the risk probability  $P(Y = 1|\mathbf{X}_k)$ , for any  $\mathbf{X}_k$ .

We assume the target model of interest  $Y|\mathbf{X}, \mathbf{B}, \mathbf{S}$  is a generalized linear model (GLM):

$$g[E(Y|\mathbf{X}, \mathbf{B}, \mathbf{S})] = \gamma_0^{S_0} + \sum_{k=1}^K \gamma_0^{S_k} S_k + \sum_{p=1}^P \gamma_{X_p}^{S_0} X_p + \sum_{k=1}^K \sum_{p=1}^P \gamma_{X_p}^{S_k} S_k X_p + \sum_{q=1}^Q \gamma_{B_q}^{S_0} B_q, \quad (4.1)$$

where  $g$  is a known link function and the terms that contain  $S_k$  indicate the changes from the internal study. We assume that the distribution of  $Y|\mathbf{X}, \mathbf{B}, \mathbf{S}$  is correctly specified, which indicates that each external population can potentially differ in intercept and  $\mathbf{X}$  covariate effect as long as those  $\mathbf{X}$ 's were used in the external model. For covariates that are unmeasured in the  $k^{\text{th}}$  external model (partial  $\mathbf{X}$  and  $\mathbf{B}$ ), we assume the covariate effects are the same as the internal study, i.e. no such population-interaction terms in model 4.1. In practice, we would force some  $\gamma_{X_p}^{S_k}$  to be zero based on prior knowledge, e.g. set  $\gamma_{X_1}^{S_k} = 0$  if we believe study  $k$  has the same  $X_1$  effect as the internal study or maybe there is not enough power to distinguish a binary  $X_1$  effect in study  $k$ . Similarly, if we believe the marginal effect of study  $k$  is the same as the internal study, we could force  $\gamma_0^{S_k}$  to be zero.

Model 4.1 is the general form of the target model as it is a saturated model allowing all intercept and possible  $\mathbf{X}$  covariates to differ across populations. A special case of model 4.1 is a logistic regression model that only allows intercept differences among populations, which represents dif-



ferent prevalence and the same covariate effect in each population:

$$\text{logit}[\Pr(Y = 1|\mathbf{X}, \mathbf{B}, \mathbf{S})] = \gamma_0^{S_0} + \sum_{k=1}^K \gamma_0^{S_k} S_k + \sum_{p=1}^P \gamma_{X_p}^{S_0} X_p + \sum_{q=1}^Q \gamma_{B_q}^{S_0} B_q. \quad (4.2)$$

We assume that the external models  $Y|\mathbf{X}_k$ 's are the best-fitted models in the class of the reduced models that was considered, but this class of reduced models may not contain the true distribution of  $Y|\mathbf{X}_k$ . One such example is when the full model is the logistic model shown in equation 4.2, but the true distribution for  $Y|\mathbf{X}_k$ 's are not logistic models as collapsibility does not hold for the logit link. We consider the fitted logistic model as the best-fitted model in the class.

### 4.2.2 Proposed Data Integration and Analysis Strategy

Figure 4.1 illustrates the proposed five-step strategy, along with the required assumptions in each step. We will first briefly introduce the steps and then expand upon the details.

- **Step 1:** Convert each external summary-level information into a set of synthetic data according to Section 2.2.2 in Chapter 2 and append each of the synthetic data sets to the internal data, from which we create a longer dataset as illustrated in Figure 4.1. The synthetic data for external study  $k$  constitutes of observed  $\mathbf{X}_k$  and the simulated value of  $Y$ . Unmeasured variables in the external populations (all  $\mathbf{B}$  and some  $\mathbf{X}$ 's) will be treated as missing data. For example, since the external study  $S=1$  used  $X_1$  and  $X_2$  to predict  $Y$ , we first replicate the observed  $(X_1, X_2)$  in  $S=0$  a large number of times (see details of the replication number in the following descriptive paragraph for step 1); we then utilize the summary information  $Y|X_1, X_2$  from external model 1 to generate the synthetic  $\hat{Y}^{S=1}$  values given  $X_1$  and  $X_2$ ; and lastly, the unmeasured variables  $X_3$  and  $\mathbf{B}$  will remain missing (Figure 4.1). Similarly for the external study  $S=2$ , we replicate the observed  $(X_1, X_3)$ , and create synthetic  $\hat{Y}^{S=2}$  values. The combined dataset is of size  $N \times (P+Q+2)$ .
- **Steps 2-3:** For the combined dataset created in step 1, multiply impute the missing covariates ignoring the outcome  $Y$  through multiple imputation by chained equation (MICE) to create  $M$  complete datasets, and then stack these  $M$  datasets to create a stacked dataset. These two steps are identical to step 1-2 of the stacked imputation approach proposed by Beesley and Taylor [2020], where the authors proposed to decompose the imputation model into two parts:  $f(X_{\text{miss}}|X_{\text{obs}}, Y) \propto f(Y|X)f(X_{\text{miss}}|X_{\text{obs}})$  while we modify it as  $f(X_{\text{miss}}|X_{\text{obs}}, \mathbf{B}, Y, S) \propto f(Y|X, \mathbf{B}, S)f(X_{\text{miss}}|X_{\text{obs}}, \mathbf{B}, S)$ . The stacked dataset is of size  $MN \times (P+Q+2)$ .



- **Step 4:** Calculate weights for each row of the stacked dataset with the weights proportional to the target model density  $f(Y|X, B, S)$ . Note that all weights need to be re-scaled to 1 within individuals in the stacked dataset. Initial parameter estimates are needed for each population as discussed in subsequent paragraphs.
- **Step 5:** Estimate the parameter  $\gamma$  of the target model 4.1 through a weighted GLM using the stacked dataset. The estimated variance of  $\gamma$  can be obtained numerically through bootstrap or analytically through several existing estimators, such as the Louis information estimator or the Jackknife estimator [Beesley and Taylor, 2020, 2021].

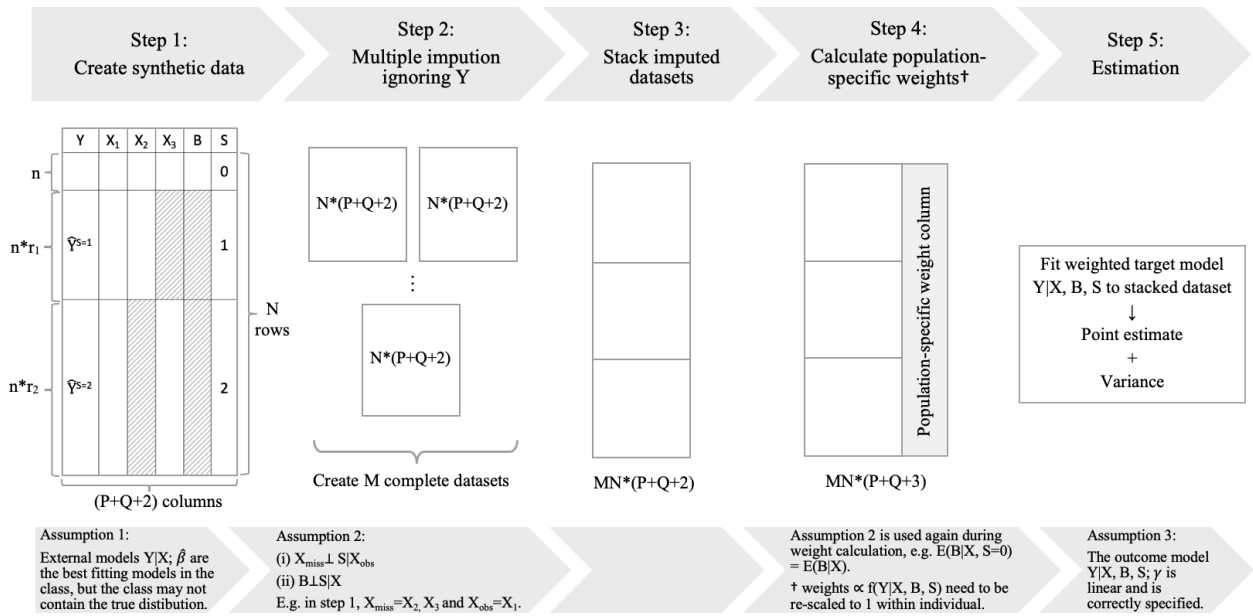


Figure 4.1: Diagram of the proposed data integration and analysis strategy.

Step 1 converts the original problem to regression modeling with missing data using a combined dataset of the internal and the synthetic data. In Chapter 2, we provided theoretical justification in special cases to show that the synthetic data method is equivalent to a constrained semi-parametric maximum likelihood approach proposed by Chatterjee et al. [2016], and it assumed the external models were the best-fitted models in the class, but the class may not contain the true distribution (Assumption 1). In finite sample size, the larger the number of replicates in each synthetic dataset (denoted as  $r_k$  in Figure 4.1), the more precision gain in the estimated coefficient of  $X$ ; when  $r_k$  goes to infinity, the precision gain by incorporating external information will converge to a constant (see details in Chapter 2). In practice, it is reasonable to set the synthetic data size, i.e.  $n*r_k$ , similar to the external study's actual study size. We will assess the performance of the proposed strategy by varying the number of replicates in the simulations described in Section 4.3.

To implement steps 2 and 3, we require some quantities to be shared across populations, since the missing covariates are completely unobserved in one population, also known as block-wise missing structure. Assumption 2 contains two parts: (i)  $\mathbf{X}_{\text{miss}} \perp\!\!\!\perp \mathbf{S} | \mathbf{X}_{\text{obs}}$ ; and (ii)  $\mathbf{B} \perp\!\!\!\perp \mathbf{S} | \mathbf{X}$ . These two assumptions imply that the conditional distribution of the missing covariates conditional on the observed covariates is the same across populations, and thus observed information can be shared across populations to impute missing covariate information. Therefore, the imputation models are  $f(\mathbf{X}_{\text{miss}} | \mathbf{X}_{\text{obs}})$  and  $f(\mathbf{B} | \mathbf{X})$  for missing  $\mathbf{X}$  and missing  $\mathbf{B}$ , respectively (e.g.  $\mathbf{X}_{\text{miss}} = [X_2, X_3]$  and  $\mathbf{X}_{\text{obs}} = X_1$  in Figure 4.1). The missing at random (MAR) assumption required by MICE is naturally satisfied as we have designed missingness, i.e., missing covariates are completely unobserved due to not being collected in the study, which is by design not related to the missing observations.

In step 4, initial parameter estimates,  $\hat{\gamma}_0$  for the internal population and  $\hat{\gamma}_k$ 's for the external populations, are needed to calculate weights that are proportional to  $f(\mathbf{Y} | \mathbf{X}, \mathbf{B}, \mathbf{S})$ . For the internal population  $\mathbf{S}=0$ , we replace  $f(\mathbf{Y} | \mathbf{X}, \mathbf{B}, \mathbf{S} = 0)$  with  $f(\mathbf{Y} | \mathbf{X}, \mathbf{B}, \mathbf{S} = 0; \hat{\gamma}_0)$ , where  $\hat{\gamma}_0$  is the internal-data-only estimates fitted on model 4.1. For external populations  $\mathbf{S}=k$ ,  $\hat{\gamma}_k$  from model  $f(\mathbf{Y} | \mathbf{X}, \mathbf{B}, \mathbf{S} = k; \hat{\gamma}_k)$  is not directly available since we only have the summary information on the reduced model  $\mathbf{Y} | \mathbf{X}_k; \hat{\beta}_k$ . As described in Section 4.2.1, the summary information will be available in the form of either parameter estimates  $\hat{\beta}_k$  (Category 1) or a risk calculator of unknown form that has the ability to estimate the probability of  $\mathbf{Y}=1$  given  $\mathbf{X}_k$  (Category 2). In the case of Category 1, we propose to derive the initial estimates  $\hat{\gamma}_k = (\hat{\gamma}_0^{S_k}, \hat{\gamma}_X^{S_k T}, \hat{\gamma}_B^{S_0 T})^T$ , where  $\hat{\gamma}_0^{S_k}$  and  $\hat{\gamma}_X^{S_k}$  are bias-corrected estimates of intercept and  $\mathbf{X}$  coefficients from  $\hat{\beta}_k$  according to Neuhaus and Jewell [1993] while  $\hat{\gamma}_B^{S_0}$  are estimated coefficients of  $\mathbf{B}$  using only the internal data. Assumption 2 is used again in this step, e.g.  $E(\mathbf{B} | \mathbf{X}, \mathbf{S} = 0) = E(\mathbf{B} | \mathbf{X})$ , so that we use the internal data to estimate the mean profile of  $\mathbf{B} | \mathbf{X}$  in the external populations (see Appendix C.1 for details). In the case of Category 2 where  $\hat{\beta}_k$  does not exist, we follow the same procedure as in Category 1 but use  $\hat{\beta}_k^{\text{synthetic}}$  instead. Specifically, we first create a large size of synthetic data  $(\hat{\mathbf{Y}}^{\mathbf{S}=k}, \mathbf{X}_k^{\text{synthetic}})$  as described in Step 1, and then we fit a GLM  $\hat{\mathbf{Y}}^{\mathbf{S}=k} | \mathbf{X}_k^{\text{synthetic}}$  with main effect using only the synthetic data and ignoring the missing data, from which we obtain  $\hat{\beta}_k^{\text{synthetic}}$ . Note that this main-effect GLM is mis-specified but the best linear model in the class. Further assessment can be found in simulation II of Section 4.3.

Step 4 extends Beesley and Taylor [2020] to allow multiple populations with potentially different covariate effects in the target model. Beesley and Taylor [2020] reduces to a special homogeneous-population case of ours when all  $\mathbf{S}$  equals to zero. In summary, we assign weights  $w_{\text{im}} = \frac{f(\mathbf{Y}_i | \mathbf{X}_{\text{im}}, \mathbf{B}_{\text{im}}, \mathbf{S}=k; \hat{\gamma}_k)}{\sum_{j=1}^M f(\mathbf{Y}_i | \mathbf{X}_{ij}, \mathbf{B}_{ij}, \mathbf{S}=k; \hat{\gamma}_k)}$  to each observation  $i$  in the  $m^{\text{th}}$  imputed dataset, where  $m \in \{1, \dots, M\}$  and  $k \in \{0, \dots, K\}$ . For internal observations,  $w_{\text{im}} = \frac{f(\mathbf{Y}_i | \mathbf{X}_{\text{im}}, \mathbf{B}_{\text{im}}, \mathbf{S}=0; \hat{\gamma}_0)}{M * f(\mathbf{Y}_i | \mathbf{X}_{\text{im}}, \mathbf{B}_{\text{im}}, \mathbf{S}=0; \hat{\gamma}_0)} = \frac{1}{M}$ .

As mentioned in Section 4.2.1, step 5 assumes that the target outcome model is linear and is correctly specified (Assumption 3). After obtaining the point estimates  $\hat{\gamma}$  by fitting the weighted

GLM, several variance estimators are available to measure the variation of  $\hat{\gamma}$ . Beesley and Taylor [2020, 2021] proposed three variance estimators including a bootstrap estimator by resampling the imputed datasets in step 2 and repeating steps 3-5. We propose a new bootstrap procedure by resampling the internal data and repeating all the steps 1-5, which is more computationally intense but empirically we find it gives more accurate estimation of the variance. We assess the performance of different variance estimators in simulations in Section 4.3.

### 4.3 Simulation Studies

Although our proposed approach can handle any target model that belongs to the class of GLM, we focus on a binary outcome and logistic regression to evaluate the performance of the proposed approach and comparison methods. In all scenarios, the internal data is of size 200, while we vary the synthetic data size from one times the internal data size (i.e.  $r_1 = r_2 = 1$  and  $N=n*[1 + r_1 + r_2]=600$ ) to 10 times (i.e.  $r_1 = r_2 = 10$  and  $N=n*[1 + r_1 + r_2]=4,200$ ) for each external population. We implement four methods, where method (1) is the benchmark, method (2) is the proposed approach, and methods (3) and (4) are two common approaches to analyze the combined dataset created in step 1 of Section 4.2.2:

1. Internal data only: fit the target model on the internal data  $S=0$  only, without incorporating external information;
2. Proposed method: we implement it through MICE in R software. For example, in Figure 1 step 1, the imputation models for  $X_2$ ,  $X_3$  and  $B$  are  $(X_2|X_1, X_3, B)$ ,  $(X_3|X_1, X_2, B)$ , and  $(B|X_1, X_2, X_3)$ , respectively. Weights are calculated after imputation;
3. FCS: imputation through fully conditionally specification (FCS) by specifying an imputation model for each missing variable conditional on all the observed covariates and the outcome  $Y$ , and iteratively generate imputed values [Van Buuren et al., 2006]. For example, in Figure 1 step 1, the imputation models for  $X_2$ ,  $X_3$  and  $B$  are  $(X_2|X_1, X_3, B, Y)$ ,  $(X_3|X_1, X_2, B, Y)$ , and  $(B|X_1, X_2, X_3, Y)$ , respectively;
4. IMB: “imputation by ordered monotone blocks (IMB)” strategy to handle block-wise missingness proposed by Li et al. [2014]. In our case, it sequentially imputes missing covariates starting with the variable with minimum missingness conditional on the observed data, outcome, and newly imputed data. We implement IMB through MICE by specifying a different imputation model compared with FCS, e.g. in Figure 1 step 1, the imputation models for  $X_2$ ,  $X_3$  and  $B$  are  $(X_2|X_1, Y)$ ,  $(X_3|X_1, X_2, Y)$ , and  $(B|X_1, X_2, X_3, Y)$ , respectively.

M=100 imputations are used for all multiple imputation. For FCS and IMB, we fit the same target model as the proposed method but without weights, and calculate the variance via Rubin's combining rules [Little and Rubin, 2002].

### 4.3.1 Simulation Settings

We provide two representative examples in Simulation I and II to illustrate how to handle the two categories of external summary-level information, respectively (Figure 4.2). Additional simulation results to assess various settings and violations of assumptions can be found in Appendix C.2.

		<b>Simulation I*:</b> Ideal case with continuous B1 and binary B2	<b>Simulation II:</b> External model 2 was derived by fitting a random forest model to a large dataset where the underlying generative model is a logistic regression model that contained quadratic and interaction terms
<b>Structure of the assumed model</b>	Internal	logit $\Pr[Y=1   X_1, X_2, B_1 \text{ (cont.)}, B_2 \text{ (binary)}]$	logit $\Pr[Y=1   X_1, X_2, X_3, X_4, B_1 \text{ (cont.)}, B_2 \text{ (binary)}]$
	External 1	logit $\Pr(Y=1   X_1; \beta_1)$	logit $\Pr(Y=1   X_1, X_2, X_3; \beta_1)$
	External 2	logit $\Pr(Y=1   X_1, X_2; \beta_2)$	A random forest model using $X_1, X_2, X_3$ , and $X_4$ to predict the probability of $Y=1$
<b>Covariate distribution</b>		$(X_1, X_2, B_1) \sim N(0, 1)$ , correlation 0.3 $B_2 \sim \text{Ber}[1 + \exp^{-1}(0.1X_1 + 0.2X_2 + 0.3B_1)]$	$(X_1, X_2, X_3) \sim N(0, 1)$ , correlation 0.3 $X_4 \sim N(0.1X_1 + 0.1X_2 + 0.1X_3, 1)$ $B_1 \sim N(0.2X_1 + 0.2X_2 + 0.2X_3 + 0.1X_4, 1)$ $B_2 \sim \text{Ber}[1 + \exp^{-1}(0.2X_1 + 0.2X_2 + 0.2X_3 + 0.1X_4 + 0.1B_1)]$
<b>Generative true outcome model</b> logit $\Pr(Y=1   X, B)$	Internal	$-1 - X_1 - X_2 - B_1 - B_2$	$-1 - X_1 - X_2 - X_3 - X_4 - B_1 - B_2$
	External 1	$1 - X_1 - X_2 - B_1 - B_2$	$2 - X_1 - X_2 - X_3 - X_4 - B_1 - B_2$
	External 2	$3 - X_1 - X_2 - B_1 - B_2$	$3 - X_1 - X_2 - X_3 - X_4 - B_1 - B_2 + 0.1(X_1^2 + X_2X_3)$
<b>Available information</b>	Internal	Individual data ( $Y, X_1, X_2, B_1, B_2$ )	Individual data ( $Y, X_1, X_2, X_3, X_4, B_1, B_2$ )
	External 1	$\hat{\beta}_1 = (0.32, -1.19)^T \dagger$	$\hat{\beta}_1 = (1.11, -1.07, -1.16, -1.06)^T \dagger$
	External 2	$\hat{\beta}_2 = (2.11, -1.13, -1.12)^T \dagger$	A random forest risk calculator that will provide the estimated probability $Y=1$ , given $X_1, X_2, X_3$ , and $X_4$ ¶
<b>Target model</b> logit $\Pr(Y=1   X, B, S)$		$\gamma_0^{S_0} + \gamma_0^{S_1} S_1 + \gamma_0^{S_2} S_2$ $+ \gamma_{X_1}^{S_0} X_1 + \gamma_{X_2}^{S_0} X_2 + \gamma_{B_1}^{S_0} B_1 + \gamma_{B_2}^{S_0} B_2$	$\gamma_0^{S_0} + \gamma_0^{S_1} S_1 + \gamma_0^{S_2} S_2$ $+ \gamma_{X_1}^{S_0} X_1 + \gamma_{X_2}^{S_0} X_2 + \gamma_{X_3}^{S_0} X_3 + \gamma_{X_4}^{S_0} X_4 + \gamma_{B_1}^{S_0} B_1 + \gamma_{B_2}^{S_0} B_2$
<b>Evaluation metrics</b>		- Absolute bias - Variance estimation - Empirical variance of point estimates	- Area under the curve (AUC) - Sum of squared error (SSE) - Scaled Brier Score (BS)

\* Additional simulation scenarios to assess the performance of the proposed approach are presented in Section 1 of Web Supplemental.  
†  $\hat{\beta}_k$ ,  $k=1$  or  $2$ , is estimated from a large dataset that follows the true data generating mechanism.  
¶ The random forest model is fitted on a large dataset of  $(Y, X_1, X_2, X_3, X_4)$  that follows the true data generating mechanism.

Figure 4.2: Simulation settings snapshot.

- **Simulation I:** Idealized case where the internal data contains  $(Y, X_1, X_2, B_1[\text{continuous}], B_2[\text{binary}])$ , and two external models have been fitted to very large datasets that is sampled from the true data generating mechanism. The external models provided parameter estimates  $\hat{\beta}_1$  and  $\hat{\beta}_2$  from logistic regression models  $Y|X_1$  and  $Y|X_1, X_2$ , respectively. In all three populations,  $X_1, X_2$  and  $B_1$  follows a standard multivariate normal distribution with zero-

mean, standard deviation 1, and 0.3 correlation, while  $B_2$  follows a Bernoulli distribution  $B_2|X_1, X_2, B_1 \sim \text{Ber}([1 + \exp^{-1}(0.1X_1 + 0.2X_2 + 0.3B_1)])$ . As shown in Figure 4.2, the three populations have similar generative outcome models with different intercepts, i.e. -1, 1 and 3, which give prevalence of  $Y=1$  of 0.3, 0.57, and 0.81, respectively. The target model is a logistic regression with different intercepts of the form  $\text{logit}[\Pr(Y = 1|\mathbf{X}, \mathbf{B}, \mathbf{S})] = \gamma_0^{S_0} + \gamma_0^{S_1}S_1 + \gamma_0^{S_2}S_2 + \gamma_{X_1}^{S_0}X_1 + \gamma_{X_2}^{S_0}X_2 + \gamma_{B_1}^{S_0}B_1 + \gamma_{B_2}^{S_0}B_2$ .

*Evaluation metrics:* We assess this simulation in terms of absolute bias, the estimated variance from bootstrap and other comparisons, and the empirical variance of point estimates.

- **Simulation II:** External model 2 was derived by fitting a random forest model to a large dataset where the underlying true generative model is a logistic regression model that contained quadratic and interaction terms. Specifically, the internal study contains complete data of  $(Y, X_1, X_2, X_3, X_4, B_1[\text{continuous}], B_2[\text{binary}])$ , and the two external models are available in different forms of summary information, external model 1 that provides  $\hat{\beta}_1$  from a logistic regression model  $Y|X_1, X_2, X_3$  and external model 2 that can provide the estimated probabilities of  $Y=1$  given  $X_1, X_2, X_3$  and  $X_4$  through a fitted random forest model. In all three populations,  $X_1, X_2$  and  $X_3$  follows a standard multivariate normal distribution with zero-mean, standard deviation 1, and 0.3 correlation, while  $X_4$  and  $B_1$  each follows a conditional normal distribution,  $X_4|X_1, X_2, X_3 \sim N(0.2\sum_{p=1}^3 X_p, 1)$  and  $B_1|\mathbf{X} \sim N(0.2\sum_{p=1}^3 X_p + 0.1X_4, 1)$ , and  $B_2$  follows a Bernoulli distribution  $B_2|\mathbf{X}, B_1 \sim \text{Ber}(\{1 + \exp^{-1}[0.2\sum_{p=1}^3 X_p + 0.1(X_4 + B_1)]\})$ , respectively. Similar to simulation I, the true generative distributions of  $Y$  in the internal and external population 1 shared the same main covariate effect but have different intercepts (-1 and 2, which corresponds to prevalence 0.3 and 0.65), while external model 2 additionally contains a quadratic term and an interaction (with intercept 3 that corresponds to prevalence 0.73). The target model is a logistic regression with the form  $\text{logit}[\Pr(Y = 1|\mathbf{X}, \mathbf{B}, \mathbf{S})] = \gamma_0^{S_0} + \gamma_0^{S_1}S_1 + \gamma_0^{S_2}S_2 + \gamma_{X_1}^{S_0}X_1 + \gamma_{X_2}^{S_0}X_2 + \gamma_{X_3}^{S_0}X_3 + \gamma_{X_4}^{S_0}X_4 + \gamma_{B_1}^{S_0}B_1 + \gamma_{B_2}^{S_0}B_2$ .

As described in step 4 of Section 4.2.2,  $\hat{\beta}_1$  can be directly used to calculate weights for  $S=1$  while we need to estimate  $\beta_2^{\text{synthetic}}$  to calculate weights for  $S=2$ . To obtain  $\beta_2^{\text{synthetic}}$ , we first generate a large synthetic data set  $(\hat{Y}^{S=2}, X_1^{\text{synthetic}}, \dots, X_4^{\text{synthetic}})$  by replicating the observed  $(X_1, X_2, X_3, X_4)$  and generating  $\hat{Y}^{S=2}$  values through the available random forest model, and then fit a main effect logistic model  $\hat{Y}^{S=2}|X_1^{\text{synthetic}}, \dots, X_4^{\text{synthetic}}$  using only the synthetic data and ignoring the missing  $B_1$  and  $B_2$ .

*Evaluation metrics:* Since prediction accuracy will be the main goal in such situation in practice, we evaluate this simulation using three prediction metrics over a validation data of size  $N_{\text{test}} = 2,000$ : Area under the curve (AUC); Sum of squared error (SSE)  $= \frac{1}{N_{\text{test}}} \sum_{i=1}^{N_{\text{test}}} (\hat{p}_i -$

$p_{i0})^2$ , where  $\hat{p}_i$  and  $p_{i0}$  denotes the estimated and true probability of  $Y_i = 1$  given  $X_i$  and  $B_i$ , respectively; and Scaled Brier Score (BS):  $= \sum_{i=1}^{N_{\text{test}}} (Y_i - \hat{p}_i)^2 / \sum_{i=1}^{N_{\text{test}}} (Y_i - \bar{Y})^2$ , where  $\bar{Y} = \frac{1}{N_{\text{test}}} \sum_{i=1}^{N_{\text{test}}} Y_i$ .

### 4.3.2 Simulation Results

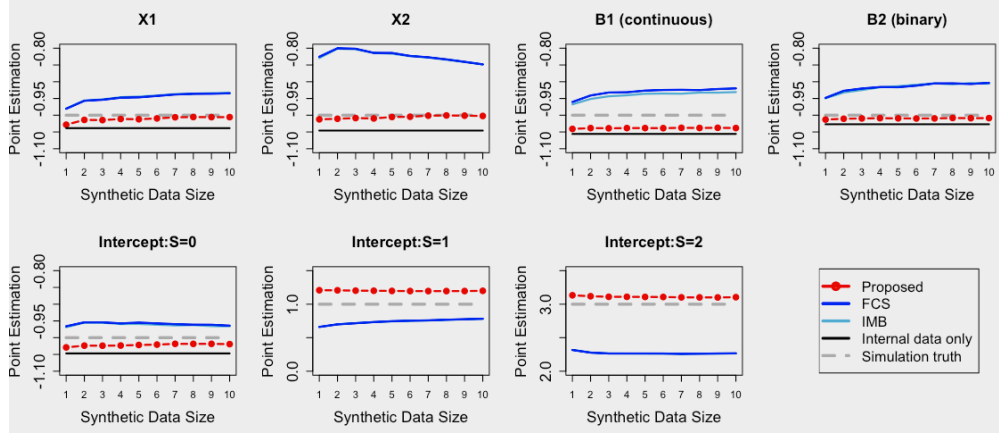
Figure 4.3 shows the average results of the target model parameter estimates across 500 simulated datasets for simulation I, including point estimates in Figure 4.3a, variance estimators versus the empirical variance of the point estimates in Figure 4.3b, and the comparison of different variance estimators for the proposed strategy in Figure 4.3c. This figure appears in color in the electronic version of this article, and color refers to that version. Figure 4.3a shows that FCS (dark blue curve) and IMB (light blue curve) have similarly biased estimates, indicating these traditional imputation strategies can not distinguish heterogeneous population effects, while the proposed method (red dotted curve) always shows close results to the truth (grey dashed curve) for all covariates, especially  $X_2$  and external intercepts  $S=2$  where other methods show severe bias. For example, the absolute bias of  $X_2$  coefficient estimates can be up to 0.2 for both FCS and IMB while it's only 0.01 for the proposed method.

As shown in Figure 4.3b, each color denotes one distinct method, along with one solid curve represents the variance estimator, and one dashed curve represents the Monte Carlo empirical variance of the point estimates. If the variance is correctly estimated, the solid curve should be approximately equal to the corresponding dashed one, which is true for all methods. As expected, all methods show precision gain in estimated X coefficients compared to the internal data only (the longer the distance to the black internal-data-only curve, the larger the precision gain) while no precision gain is found in B covariates and the intercept due to no external added information and allowing population-specific effects, respectively. The proposed method has over 50% efficiency gain in estimated X coefficients compared to the internal data. We see FCS and IMB have larger precision gain in both estimated X coefficients than the red proposed curve, which may be explained by bias-variance trade-off as they also have larger bias in the corresponding point estimates in Figure 4.3a. We will discuss the underlying statistical reason in Section 4.5.

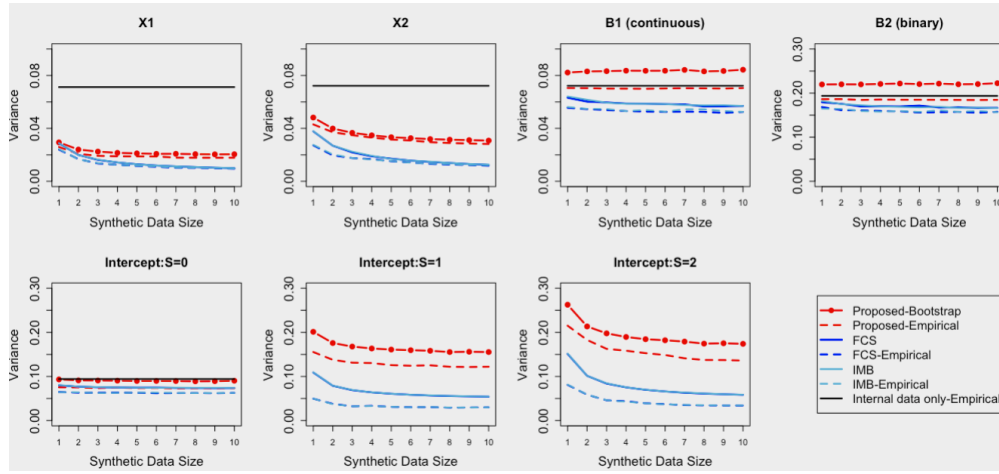
As shown in Figure 4.3c, we show the result of the Louis information estimator (StackImpute-Louis), one of the three variance estimators proposed by Beesley and Taylor [2020] as they always have similar performances. We have shown in Chapter 2 that when the synthetic data size goes to infinity, the precision gain we achieve in X covariates will converge to a constant, which is shown as the gradually stable trend of the grey curve (Monte Carlo empirical variance of the point estimates, also serves as the empirical truth here). When the synthetic data size increases from one times the internal data size to 10 times for each external study (i.e. total missing rate increases from 66.6%

to 95%), the StackImpute variance estimator and Rubin's rule variance continuously underestimate the empirical truth. On the contrary, the proposed variance estimator by bootstrapping the whole proposed procedure is always close to the empirical truth, especially in X covariates where the bias in other methods can be 10 times higher than the proposed method (i.e. 0.02 versus 0.002 in absolute bias) and could be even larger when the synthetic data size keep increasing. Moreover, in estimating the variation of the coefficient corresponding to  $B_2$ , the proposed method has stable bias around 0.035 while the other methods exhibit substantially larger bias. Note that the internal-data-only results (black solid curve) does not exist in external intercepts as they were not available in the internal data, whereas the bias of the internal data estimates is due to the small sample bias compared with the simulation truth.

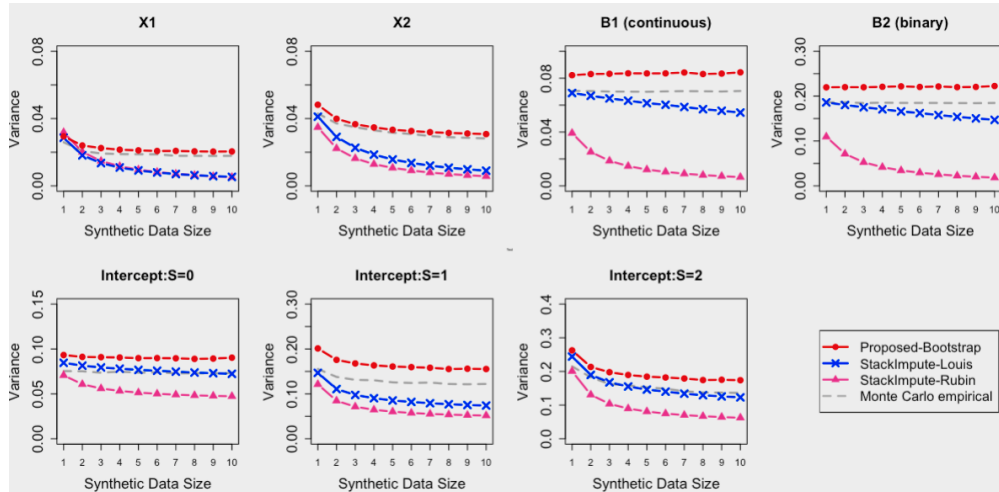




(a) Point estimates



(b) Variance estimators vs. the empirical variance



(c) Different variance estimators of the proposed method

Figure 4.3: Visualization of simulation I results over increasing synthetic data size (a) point estimates (b) variance estimators vs. the empirical variance (c) different variance estimators of the proposed strategy.



Figure 4.4 shows the performance of each method in Simulation II over increasing synthetic data size on a validation dataset that follows the true data generating mechanism, with three prediction metrics on the row and three populations on the column. In general, the results in Figure 4.4 are in line with Figure 4.3a, implying that the proposed method has better overall prediction performance compared with others. Specifically, in the first column (internal population  $S=0$ ), all methods incorporating external information have consistently better prediction ability (larger AUC, smaller SSE and smaller BS) than using the internal data only (the dashed grey line). While all methods have similar performance in terms of AUC (first row) and predicting internal population (first column), the proposed approach outperforms others in terms of SSE and BS in predicting external populations, especially external study 2 where the true parameter values are quite different from the internal study values (the proposed method has up to 41% more improvement in SSE and 19% more improvement in BS compared with FCS and IMB). The proposed method shows a modest improvement in performance as the size of the synthetic data increases, e.g. for  $S=2$ , SSE decreases 12.9% from 9.3 to 8.1 in the proposed method when the synthetic data size increases from one to 10. In this particular scenario, there is little gain in performance with synthetic data more than 4 times the internal data size. Note that it is hard to distinguish FCS and IMB in the figure as they have very close results.

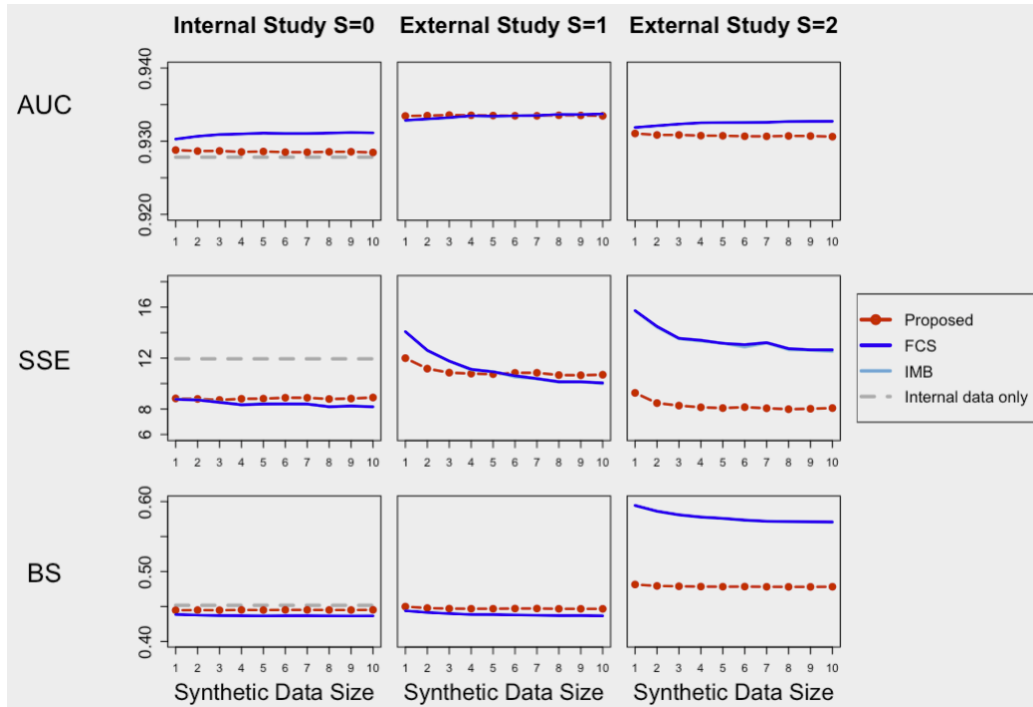


Figure 4.4: Visualization of prediction metrics over increasing synthetic data size for simulation II. Larger AUC, smaller SSE and smaller BS represents better prediction.

## 4.4 Application to Prostate Cancer Data

We apply the proposed method to predict the risk of high-grade prostate cancer (Gleason score over 6) using an internal dataset containing patients from three United States academic institutions [Tomlins et al., 2015] and two external risk calculators, one Prostate Cancer Prevention Trial risk calculator established from a United States population [Thompson et al., 2006] and another European Randomized Study of Screening for Prostate Cancer risk calculator 3 established from a European population [Roobol et al., 2012].

The external risk calculators each used slightly different predictors to predict the same outcome through logistic regression, which we will denote as PCPThg and ERSPC, respectively. PCPThg and ERSPC both used prostate-specific antigen level (PSA) and digital rectal examination findings (DRE) as one of the predictors, and PCPThg also used age, race (African American or not) and prior biopsy results while ERSPC additionally used transrectal ultrasound prostate volume (TRUS-PV):

- PCPThg:  $\text{logit}(p_i) = -3.69 + 0.89\log_2(\text{PSA}_i) + \text{DRE}_i + 0.03\text{Age}_i + 0.96\text{Race}_i - 0.36\text{Biopsy}_i$ ;
- ERSPC:  $\text{logit}(p_i) = -3.16 + 1.18\log_2(\text{PSA}_i) + 1.81\text{DRE}_i - 1.51\log_2(\text{TRUS-PV}_i)$ ,

where  $p_i$  is the probability of observing high-grade prostate cancer for subject  $i$ . A more detailed description of the above two equations can be found in Appendix B.5 of Chapter 3.

In the internal individual-level data, in addition to all the predictors used in PCPThg and ERSPC, we also have data on two new biomarkers, prostate cancer antigen 3 (PCA3) and TM-PRSS2:ERG (T2:ERG) gene fusions, that were prognostic of prostate cancer [Tomlins et al., 2015, Truong et al., 2013]. Therefore, in Table 4.1, we present the results of two target models using a total of eight predictors, including these two new biomarkers, one model only allows the intercept to be different across populations (“different intercept only” model corresponds to model 4.2), and another flexible model allows all possible covariates that used in the external populations to have population-specific effects (“different intercept and covariates” model that corresponds to model 4.1). The detailed model forms can be found in the table legend. A total of 678 male patients who had complete data were included in the internal data set provided by Tomlins et al. [2015], and an additional 1,174 patients’ data were independently collected from seven community clinics throughout the United States for validation.

The grey empty blocks in Table 4.1 imply that these predictors were not used in the certain external model, and thus in the proposed method we assume they have the same coefficient as the internal population (grey blocks with values). Results in Table 4.1 show (i) we will gain precision of the estimated coefficients by incorporating external model information (smaller SE highlighted in green), e.g. in the different-intercept-only model, the bootstrap SE of  $\log_2(\text{PSA})$  reduces from

Table 4.1: Results of the data example for predicting the risk of high-grade prostate cancer. Internal dataset of size  $n=678$ ; validation dataset of size  $N_{\text{test}}=1,174$ ; SE, standard error; green represents good performance with large precision gain or better model calibration, grey represents certain predictor was not used in the external calculator, yellow represents population-specific effect different from the internal population, and red represents poor performance compared with direct regression; SE in the proposed method are (StackImpute SE) [bootstrap SE from 500 replicates].

	PCPThg		ERSPC		Proposed method			
	Original		Estimated		Direct regression		Different intercept only†	
	Estimate	SE	Estimate	SE	Estimate	SE	Internal	ERSPC
Intercept	-3.686	(.115)	-1.409	(.115)	-3.16	(.111)	-4.157	-6.884
$\log_2(\text{PSA})$	0.894	(.124)	0.735	(.124)	1.176	(.124)	1.082	1.082
DRE	1	(.257)	1.145	(.257)	1.813	(.269)	1.118	1.118
Age	0.03	(.012)	0.033	(.012)	/	(.012)	0.030	0.030
Biopsy	-0.36	(.272)	-1.444	(.272)	/	(.272)	-0.618	-0.618
Race	0.96	(.288)	0.442	(.288)	/	(.288)	0.434	0.434
$\log_2(\text{TRUS-PV})$	/	/	/	/	-1.514	(.225)	-1.553	-1.553
$\log_2(\text{PCA3+1})$	/	/	/	/	/	(.225)	-1.553	-1.553
$\log_2(\text{T2:ERG+1})$	/	/	/	/	/	(.225)	-1.553	-1.553
PredictionAUC	0.700	0.707	0.707	0.723	0.720	0.723	0.804	0.804
metrics** Scaled Brier score	1.059	0.987	0.987	0.994	0.954	0.994	0.915	0.915
$\gamma_0^{S_0} + \gamma_0^{S_{\text{PCPThg}}}$	$\gamma_0^{S_{\text{PCPThg}}} + \gamma_0^{S_{\text{ERSPC}}}$	$\gamma_0^{S_{\text{PCPThg}}} + \gamma_0^{S_{\text{ERSPC}}}$	$\gamma_0^{S_{\text{PCPThg}}} + \gamma_0^{S_{\text{ERSPC}}}$	$\gamma_0^{S_{\text{PCPThg}}} + \gamma_0^{S_{\text{ERSPC}}}$	$\gamma_0^{S_{\text{PCPThg}}} + \gamma_0^{S_{\text{ERSPC}}}$	$\gamma_0^{S_{\text{PCPThg}}} + \gamma_0^{S_{\text{ERSPC}}}$	$\gamma_0^{S_{\text{PCPThg}}} + \gamma_0^{S_{\text{ERSPC}}}$	$\gamma_0^{S_{\text{PCPThg}}} + \gamma_0^{S_{\text{ERSPC}}}$
$\gamma_1^{(*)} + \gamma_1^{S_{\text{PCPThg}}}$	$\gamma_1^{(*)} + \gamma_1^{S_{\text{PCPThg}}}$	$\gamma_1^{(*)} + \gamma_1^{S_{\text{PCPThg}}}$	$\gamma_1^{(*)} + \gamma_1^{S_{\text{PCPThg}}}$	$\gamma_1^{(*)} + \gamma_1^{S_{\text{PCPThg}}}$	$\gamma_1^{(*)} + \gamma_1^{S_{\text{PCPThg}}}$	$\gamma_1^{(*)} + \gamma_1^{S_{\text{PCPThg}}}$	$\gamma_1^{(*)} + \gamma_1^{S_{\text{PCPThg}}}$	$\gamma_1^{(*)} + \gamma_1^{S_{\text{PCPThg}}}$
**prediction metrics on validation dataset								

0.146 to 0.070 while it reduces for DRE from 0.299 to 0.139, compared with direct regression; (ii) when allowing population-specific effects (yellow blocks), we will not expect to see much precision gain due to variance-bias trade-off, e.g. both precision gain of  $\log_2(\text{PSA})$  and DRE we see in the different-intercept-only model diminishes in the different-intercept-and-covariates model; and (iii) similar to the results in simulations, the analytical SE (in the round parenthesis) tends to provide a smaller estimation than the bootstrap SE (in the squared parenthesis) and potentially underestimates the true variability of the estimates.

In the prediction metrics row, we show AUC (higher value represents better discrimination) and scaled Brier score (smaller value means better calibration) calculated using the validation cohort, where the red blocks implies slightly worse overall predictive performance of PCPThg population compared to direct regression, e.g. 0.9% reduction of AUC. This may be because that the validation cohort may represent a moderately different population than the training cohort as it has different baseline distribution as noted by Tomlins et al. [2015], which may also explain why the fitted model for the European population has better performance on the validation data.

## 4.5 Discussion

**Flexibility in external models and populations:** The proposed approach adds to the existing research on integrating external summary information into the internal study. It can develop improved models and provide statistical inference not only for the internal population but also for the external populations. The parameters of the external population models are allowed to differ from those of the internal population. This new strategy has the appealing feature of being able to make use of external information that comes in the form of a “black box” algorithm, i.e. an algorithm that provides a predicted probability, but the underlying model is not necessarily simple or transparent or even known. The key aspect is that the external information allows the creation of synthetic data. We further summarize some key points and concluding remarks of the proposed strategy in the subsequent paragraphs.

**Using partial information through data integration:** The proposed method can integrate summary information from multiple external models that each uses different covariates  $\mathbf{X}_k \in \mathbf{X}$  into the current study. The simulation and real data analysis showed expected results that we can only gain precision on the estimated coefficient of  $\mathbf{X}$  but not the  $\mathbf{B}$  coefficients that are only available internally, even when  $\mathbf{X}$  and  $\mathbf{B}$  are correlated. This is consistent with the theoretical results in Dai et al. [2012] that the MLE estimator  $\hat{\beta}$  and  $\hat{\gamma}_B$  are always asymptotically independent under regularity conditions, where  $\hat{\beta}$  is the estimates of intercept and  $\mathbf{X}$  coefficients in model  $Y|\mathbf{X}; \beta$  and  $\hat{\gamma}_B$  is the estimated coefficient of  $\mathbf{B}$  from model  $Y|\mathbf{X}, \mathbf{B}; \gamma$ , respectively.

**Principled inference post-imputation:** Since the proposed strategy uses the stacked multiple

imputation (StackImpute) proposed by Beesley and Taylor [2020], it also borrows strengths from StackImpute to avoid incompatibility between the imputation model and the analysis model, and can accommodate complicated outcomes such as the time-to-event outcome in survival models. In the proposed method, we introduced two types of variance estimation, the proposed bootstrap variance and the analytical variance estimators proposed in Beesley and Taylor [2020, 2021]. In general, the bootstrap variance can provide more valid variance estimation but may be more computationally intense compared with others, while analytical estimates are fast to compute but may be biased. Based on simulation results, when the predictors have small covariate effects (simulation in Appendix C.2.2), the bias of the analytical variance estimates is small.

**Improve not just current study but external model predictions:** To our knowledge, very few existing approaches can allow different population effects through regression analysis in data fusion, let alone improving external model predictions. It is worth noting that this same problem setting has a wide range of applications beyond regression analysis. Several approaches are proposed in the causal inference field to estimate the average causal effect, aiming to incorporate the supplementary information from the validation dataset to the main dataset and allowing heterogeneous treatment effects among different data sources. For example, Antonelli et al. [2017] proposed a unified Bayesian imputation framework built upon the work by Wang et al. [2012b], introducing a dependence parameter to represent the prior odds of including a predictor in the outcome model given that it is in the exposure model, and assuming different population indicators in the outcome model. Yang and Ding [2020b] posited a stochastic framework on the estimators from different data sources which flexibly leverages the supplementary information from validation datasets and they used a sensitivity parameter to quantify the systematic difference among data sources. Similarly, Huang and Qin [2020] and Chen et al. [2020] addressed this same problem in the implementation of survival data without assuming comparability among data sources.

**Allow for violation of transportability assumption:** The transportability assumption is common in data integration and causal inference when certain variables are not mutually observed across populations [Rassler, 2004, Reiter, 2012, Bareinboim and Pearl, 2013]. Compared with more strict transportability assumed in the literature, e.g. Chatterjee et al. [2016] assumed transportability of the joint distribution of  $Y, X, B$  while Antonelli et al. [2017] assumed conditional transportability of  $X_{\text{miss}}|X_{\text{obs}}, Y$ , we only require conditional transportability among covariates (Assumption 2). While the simulation results (Appendix C.2.4) suggest that violating this assumption could have a mild impact, one can consider applying additional shrinkage methods such as the empirical Bayes approaches proposed by Estes et al. [2017] and us (in Chapter 3) after obtaining estimates from the proposed approach, which can empirically strike a balance between bias and efficiency when the transportability between populations is unclear.

**Limitations:** While the results of the simulation study suggest that the proposed strategy has

promising performance in providing both accurate statistical inference and prediction compared with comparison methods, some limitations are worth noticing. Particularly, the proposed method relies on good initial estimates for each external population. We propose to use a geometric approach by utilizing the observed data relationship to map the parameter estimates in the reduced model (i.e.  $Y|X; \beta$ ) to the target model (i.e.  $Y|X, B; \gamma$ ). While the simulation results show promising performance, caution must be exercised during implementation when the true underlying relationship is hard to verify. In a special case where the external study has the same population distribution as the internal population, the internal data estimates can directly serve as the initial estimates for the external population.

**Future directions:** An interesting extension of the proposed method is to accommodate the situation where selection bias exists and selection probability or survey weights are available for each observation in the internal population. In theory, the proposed method can be adapted to accommodate this by replacing the synthetic  $Y$  values with the inverse probability-weighted or survey weights-weighted synthetic  $Y$  values in step 1 of the proposed strategy. Alternatively, instead of copying the whole internal  $X$ 's multiple times to create the same  $X$  distribution as the internal population, one can consider proportionally creating synthetic  $X$  through the given weights to recovery the representative distribution in the external populations. Further investigation is needed to evaluate this. Furthermore, if the exposure indicator is available as a covariate in all populations, one can also use the regression estimates from the proposed method to calculate the estimated average causal effect by averaging over the joint distribution of  $(X, B)$ . On the contrary, it is unclear whether we can directly use the intermediate parameter estimates from causal inference methods for regression inference. For example, the causal inference approach-guided Bayesian method adjusting for unmeasured confounding [Antonelli et al., 2017] aimed at obtaining unbiased causal effects by averaging selective regression models, can also produce the regression estimates for all covariates, which is the same as our goal. We attempted to modify their method and code to serve our purpose but the results did not seem promising, this may have been because a direct comparison of the performance of two approaches is not appropriate when they have different goals.

## 4.6 Software and Publication

R package SynDI implementing the proposed method can be found on GitHub at <https://github.com/umich-biostatistics/SynDI>. The content of this chapter has been submitted for publication. A preprint is available on arXiv at <https://arxiv.org/abs/2106.06835>.



## CHAPTER 5

### Discussion

In this dissertation, we propose three statistical methods in Chapter 2, 3 and 4, respectively, to incorporate external summary-level information into the regression analysis of a current study. In Chapter 2, we propose a novel synthetic data method that can convert the external model information into synthetic data, where the only required external model information is the ability to generate predictions of the outcome given covariates. We introduce the synthetic data method under the setting of one single external study. In Chapter 3, we propose a meta-inference framework using an empirical Bayes (EB) approach considering multiple external studies, where each of the external models provides a set of regression estimates. The meta-inference framework is efficient and robust to make valid inference of the internal study as it selects the compatible external model estimates for the internal study and adaptively assigns weights according to compatibility. In Chapter 4, we extend the synthetic data method to accommodate the situation with multiple external models and further allow for heterogeneous covariate effects across external populations.

**Different approaches to solve the same genre of problem:** The chapters of this dissertation form a coherent whole, with the consistent goal to solve the same data integration problem but each has different priorities. All three proposed methods are flexible to incorporate external models that use a slightly different set of covariates, as long as the covariates are the subset of the observed covariates in the internal study. The synthetic data method proposed in Chapter 2 relaxes the requirements of the information that is available externally compared with the traditional constrained maximum likelihood approaches. In Chapter 3, the meta-inference framework focuses on making robust inference for the internal study, i.e. selecting the most comparable information to the internal study when multiple external models are available for use, aiming to strike a balance between efficiency gain and making valid statistical inference for the internal study. In Chapter 4, in addition to the features we captured in the previous two chapters (i.e. leveraging auxiliary information from a broad class of externally fitted model or established risk calculators of unknown form as in Chapter 2; and accommodating multiple external model information in Chapter 3), we further allow heterogeneous covariate effects across the external populations by extending the synthetic data method proposed in Chapter 2.

In summary, (i) Chapter 2 considers one single external study while Chapter 3 and 4 consider the situation where summary-level information from multiple external studies are available for use; (ii) Chapter 2 and 3 focus on making improved inference for the internal study while Chapter 4 can also make inference for the external populations; (iii) Chapter 2 and 4 can handle more flexible external information, from parametric regression model estimates to any machine learning models of unknown form, compared with Chapter 4, where regression model estimates are needed; and (iv) due to EB estimator's robust feature from shrinkage effect, Chapter 3 has less requirement on similarities across populations (i.e. collapsibility and transportability assumptions) compared to others, which we will discuss in details in the subsequent paragraph.

**Transportability and collapsibility assumptions in data integration:** The issues of transportability of distributions and collapsibility of prediction models are critical and inevitable in data integration. Transportability concerns the distributional similarity across the populations while collapsibility relates to the models being used within the population. In this dissertation, collapsibility is concerned with whether the distribution implied by the  $Y|X_k$  model is compatible with the distribution implied by the  $Y|X, B$  model. In Chapter 2 and 4, the collapsibility assumption asks for that the external models  $Y|X_k$ 's were the best-fitted models in the reduced class that was considered, but this class of reduced models might not contain the true distribution of  $Y|X_k$ . Full or partial transportability assumption (i.e. transportability of the joint distribution  $Y, X, B$  or conditional transportability of  $Y|X, B$  distribution) is commonly assumed in the literature when certain variables are not mutually observed across populations [Antonelli et al., 2017, Chatterjee et al., 2016, Estes et al., 2017] yet hard to verify in practice. In Chapter 4, we require mild conditional transportability among covariates  $X$  and  $B$  for valid imputation, violating which had a mild impact as shown in simulations. No specific transportability or collapsibility assumption is needed in Chapter 3 as the EB estimator's shrinkage effect limits the impact of external model information that is not compatible with the internal data and thus protects against the severe bias. Another advantage of the proposed methods is that in all three methods, we do not require the marginal  $X$  distribution to be the same across populations, as is assumed in other literature [Chatterjee et al., 2016, Kundu et al., 2019].

In practice, one can consider applying additional shrinkage methods such as the EB approaches proposed by Estes et al. [2017] or in Chapter 3 after obtaining estimates from other proposed approaches, which can empirically strike a balance between bias and efficiency when the transportability between populations is unclear. However, as we saw in simulations, the efficiency will decrease in order to trade for some level of robustness.

**Gaining efficiency by using partial information through data integration:** The proposed methods can integrate summary information from single or multiple external models that each uses different covariates  $X_k \in X$  into the current study. The simulations and real data analysis



in all three chapters show that we can gain around 10% to over 50% efficiency on the estimated coefficient of  $\mathbf{X}$  compared with the internal-data-only analysis, where the factors that limit the efficiency gain include that too few external models used certain covariate, the internal sample size is large so that the largest possible improvement is small, or the shrinkage method is applied to protect the potential bias. When we evaluate the prediction metrics such as SSE, the improvement compared with others is around 16%–40% while the magnitude of improvement in AUC or Brier score is relatively small around 1%–20%. We also notice an expected result that we can only gain precision on the estimated coefficient of  $\mathbf{X}$  but not the  $\mathbf{B}$  coefficients that are only available internally, even when  $\mathbf{X}$  and  $\mathbf{B}$  are correlated. This is consistent with the theoretical results in Dai et al. [2012] that the MLE estimator  $\hat{\beta}$  and  $\hat{\gamma}_B$  are always asymptotically independent under regularity conditions, where  $\hat{\beta}$  is the estimates of intercept and  $\mathbf{X}$  coefficients in model  $Y|\mathbf{X}; \beta$  and  $\hat{\gamma}_B$  is the estimated coefficient of  $\mathbf{B}$  from model  $Y|\mathbf{X}, \mathbf{B}; \gamma$ , respectively.

**Novelty of the synthetic data method:** It is worthy to highlight the advantages of the synthetic data method we proposed in Chapter 2. By creating large pseudo-data that is compatible with the externally established model, the synthetic data method naturally incorporates the external summary-level information into the internal data. Compared with converting the external information into constraint as in CML approaches, the synthetic data method not only simplified the task from solving complex constrained optimization, but also provides a potentially more flexible and general framework to handle this problem. Any tools that can handle missing data, including the multiple imputation technique we used, can be applied to analyze the combined dataset of the internal and the synthetic data. The only requirement for the synthetic data approach is the ability to generate outcome values given covariates from the information of the external models, without the need to know the exact form of the model. The extension to multiple external populations in Chapter 4 further allows researchers to make statistical inference on both the internal and the external populations. It is broadly applicable for general data types for the outcome and covariates, when covariates are of multi-dimension, and when each of the external models uses a different subset of covariates. These features are particularly appealing when both the parametric regression modeling and the machine learning algorithms with non-trivial form have been widely used in risk prediction modeling.

**Size of the internal and external studies:** In this dissertation, we assume the external models are well-established on large data that produce credible estimates or predicted outcomes. In the synthetic data methods (Chapter 2 and 4), the theoretical results suggest creating a very large size of synthetic data for each external study in order to gain the maximum possible efficiency gain on the condition that the external model  $Y|\mathbf{X}_k$  is compatible with the target model  $Y|\mathbf{X}, \mathbf{B}$ . In practice when the exact compatibility is hard to satisfy, we recommend limiting the synthetic data size similar to the external study’s actual study size. In Chapter 3, we provide an option to account

for the external model uncertainty by taking into account the variance-covariance matrix of the external model estimates.

**Design of the internal and external studies:** We have considered various sources of information variation across the models and populations in this dissertation, e.g. different forms of external model information, different subsets of covariates used by each of the external models, and heterogeneous covariate effects across populations. Caution must be exercised as there could still be fundamental variability due to design and sampling differences across studies. Examples of these include case-control studies, outcome-dependent sampling, selection bias in who is included in the dataset (either internal or external) and measurement bias during data collection, ignoring which could lead to biased inference in data integration.

**Future directions:** One possible extension of the proposed method is the application in causal inference. In all three chapters, regression estimates of the target model are our ultimate product. If the treatment indicator was available as one of the  $\mathbf{X}$  covariates, we could directly calculate the estimated average causal effect by averaging over the joint distribution of  $(\mathbf{X}, \mathbf{B})$  through the formula  $E_{\mathbf{X}, \mathbf{B}}[E(Y|\text{treatment} = 1, \mathbf{X}, \mathbf{B}) - E(Y|\text{treatment} = 0, \mathbf{X}, \mathbf{B})]$  using the regression estimates obtained from the proposed methods.

Another interesting extension is to modify the proposed methods to accommodate electronic health record (EHR) data under a slightly different setting and goal. Instead of having a moderately-sized unbiased internal dataset and several external models fitted on large data as in the current setting, we would then consider having a massively large number of internal EHR data, which could provide us potentially biased results due to selection bias, a common issue in EHR data. The goal would become to correct the bias from EHR data by incorporating the external information from models that are fitted on smaller sample sizes (compared to the large size of EHR data).

Another point of future consideration is to derive the asymptotic variance estimates of the target point estimates in Chapter 4 (i.e. the point estimates from the weighted GLM using the stacked dataset after multiple imputation). We showed in simulations that the existing fast-to-calculate variance estimators [Beesley and Taylor, 2020, 2021] only had unbiased results when the covariate effect is small or the synthetic data size is small (e.g. one times the internal data size), which contradicts the fact that larger synthetic data size can guarantee the efficiency gain. Future investigation studying the impact of the proposed weights on both the between-imputation and within-imputation variation may help researchers better understand the inference in the stacked imputation and possibly provide faster tools to estimate the variation.

In this dissertation, we consider the problem under the setting that we have a moderate-dimensional internal data, and each of the external models uses a subset of covariates in the internal study. Although the proposed methods are applicable in high-dimensional settings in theory

and we evaluated them in the settings of multi-dimension  $\mathbf{X}$  and  $\mathbf{B}$ , e.g.  $\dim(\mathbf{X})=9$  and  $\dim(\mathbf{B})=5$ , the performance of applying our methods to higher dimensional data such as the genetic data is unknown. Additional variable selection procedure could be considered if the dimension of  $\mathbf{B}$  is large.

## APPENDIX A

### Appendix of Chapter 2

#### A.1 Derivation of asymptotic variances for special case 1: Gaussian

##### A.1.1 Standard MLE:

The log-likelihood is:

$$l = l(\gamma, \sigma_\gamma^2) = \sum_{i=1}^n \log f(Y_i | X_i, B_i; \gamma, \sigma_\gamma^2) = -\frac{n}{2} \log(\sigma_\gamma^2) - \frac{1}{2\sigma_\gamma^2} \sum_{i=1}^n (Y_i - \gamma_X X_i - \gamma_B B_i)^2$$

Igorning the  $\sigma_\gamma$  row and column in the information matrix, the expected information matrix for  $\hat{\gamma}$  is:

$$\mathbf{I} = -E_{XB} \left( \frac{\partial^2 l}{\partial \gamma \gamma^T} \right) = n \frac{\sigma_X^2}{\sigma_\gamma^2} \begin{pmatrix} 1 & \theta \\ \theta & \theta^2 + \frac{\sigma_\theta^2}{\sigma_X^2} \end{pmatrix} \quad (\text{A.1})$$

Then we can obtain the asymptotic covariance matrix of  $\hat{\gamma}$  by taking the inverse of  $\mathbf{I}$ :

$$\text{Var}(\hat{\gamma}) = \mathbf{I}^{-1} = \frac{1}{n} \frac{\sigma_\gamma^2}{\sigma_\theta^2} \begin{pmatrix} \theta^2 + \frac{\sigma_\theta^2}{\sigma_X^2} & -\theta \\ -\theta & 1 \end{pmatrix}. \quad (\text{A.2})$$

Thus, the asymptotic variance of  $\hat{\gamma}_X$  and  $\hat{\gamma}_B$  are equal to  $\frac{\sigma_\gamma^2}{n\sigma_\theta^2}(\theta^2 + \frac{\sigma_\theta^2}{\sigma_X^2})$  and  $\frac{\sigma_\gamma^2}{n\sigma_\theta^2}$ , respectively.

##### A.1.2 Approach 1: Synthetic data method

If the synthetic data approach is applied, and under the assumption that the true value of  $\beta$  and  $\sigma_\beta$  are used to generate the synthetic data, then the combined data will have the same distribution as a dataset of size  $n+m$  in which  $m$  values of  $B$  have been removed. For this particular data structure, it is possible to obtain formulas for the asymptotic variance of the MLE of  $\gamma$ . In particular, Gourrier-

oux and Monfort [1981] gave the exact expression of the MLE and the corresponding asymptotic covariance in such case. The likelihood for the combined data is  $\prod_{i=1}^n f(Y_i, B_i|X_i) \times \prod_{i=n+1}^{n+m} f(Y_i|X_i)$ , which can be rewritten as  $\prod_{i=1}^{n+m} f(Y_i|X_i) \times \prod_{i=1}^n f(B_i|X_i, Y_i)$ . Based on this likelihood, they introduced a set of transformed parameters, and re-parameterized the distributions 2.5–2.7. They then identified the 1-to-1 relationship among the original parameters and the new set of parameters, which we will explain in the subsequent paragraph.

We obtain the estimators of the original parameters by the re-parameterization method, and then apply the delta method to get the asymptotic variance of  $\hat{\gamma}_B$  and  $\hat{\gamma}_X$ . According to Gouriou and Monfort [1981], we introduce a set of transformed parameters  $a, b, c, d$ , and  $e$ , and re-parameterize the distributions 2.5–2.7 as follows:

$$\begin{aligned} Y|X &\sim N(bX, a^2) \\ B|Y, X &\sim N(dY + eX, c^2) \end{aligned}$$

Then we identify the 1-to-1 relationship among the original parameters and the new set of parameters:

$$\begin{aligned} a^2 &= \sigma_\gamma^2 + \gamma_B^2 \sigma_\theta^2 \\ b &= \gamma_X + \theta \gamma_B \\ c^2 &= \frac{\sigma_\gamma^2 \sigma_\theta^2}{a^2} \\ d &= \frac{\gamma_B \sigma_\theta^2}{a^2} \\ e &= \theta - db \end{aligned} \tag{A.3}$$

The MLE of  $a, b$  and their asymptotic variances are easy to obtain from the linear model  $Y_i = bX_i + u_i$ ,  $\text{Var}(u_i) = a^2$ , where  $i = 1, \dots, n+m$ . Similarly, the MLE of  $c, d$ , and  $e$  and their asymptotic variances are easy to obtain from the linear model  $B_i = dY_i + eX_i + v_i$ ,  $\text{Var}(v_i) = c^2$  where  $i = 1, \dots, n$ . The estimators of the original parameters are obtained through the relationship derived from equations (A.3), where

$$\begin{aligned} \theta &= bd + e \\ \sigma_\theta^2 &= a^2 d^2 + c^2 \\ \gamma_B &= \frac{a^2 d}{\sigma_\theta^2} \\ \gamma_X &= b - \gamma_B \theta \\ \sigma_\gamma^2 &= \frac{a^2 c^2}{\sigma_\theta^2} \end{aligned}$$

and the asymptotic variance of  $\hat{\gamma}_X$  and  $\hat{\gamma}_B$  can be derived using the delta method:

$$\begin{cases} \text{Var}(\hat{\gamma}_B) = \frac{1}{n} \left[ \frac{\sigma_\gamma^2}{\sigma_\theta^2} + 2(\lambda - 1) \frac{\gamma_B^2 \sigma_\gamma^4}{\sigma_\beta^4} \right] \\ \text{Var}(\hat{\gamma}_X) = \theta^2 \text{Var}(\hat{\gamma}_B) + \frac{1}{n} \frac{\sigma_\gamma^2}{\sigma_X^2} \frac{\lambda \sigma_\gamma^2 + \gamma_B^2 \sigma_\theta^2}{\sigma_\beta^2}, \end{cases}$$

Therefore, we find the relative efficiency gain of  $\text{Var}(\hat{\gamma}) = \text{Var}(\hat{\gamma}_X, \hat{\gamma}_B)^T$  by adding  $m$  synthetic data observations compared to the original dataset of size  $n$  is

$$\text{ARE}[\text{Var}(\hat{\gamma})] = \mathbf{1} - (1 - \lambda) \begin{pmatrix} \frac{2\sigma_\gamma^2 \gamma_B^2 \sigma_\theta^2}{\sigma_\beta^4} + \frac{\sigma_\theta^2}{\sigma_\theta^2 + \sigma_X^2 \theta^2} \frac{\sigma_\gamma^2 (2\sigma_\gamma^2 - \sigma_\beta^2)}{\sigma_\beta^4} \\ \frac{2\gamma_B^2 \sigma_\theta^2 \sigma_\gamma^2}{\sigma_\beta^4} \end{pmatrix},$$

where  $\theta = \frac{\beta - \gamma_X}{\gamma_B}$ , and  $\sigma_\theta^2 = \frac{\sigma_\beta^2 - \sigma_\gamma^2}{\gamma_B^2}$ . When  $m$  gets very large such that  $\lambda \approx 0$ ,  $\text{ARE}[\text{Var}(\hat{\gamma}_X)] = 1 - \frac{\sigma_\theta^2}{\sigma_\theta^2 + \sigma_X^2 \theta^2} \frac{\sigma_\gamma^2}{\sigma_\gamma^2 + \gamma_B^2 \sigma_\theta^2} \frac{2\sigma_\gamma^2 - \sigma_\beta^2}{\sigma_\beta^2} - \frac{2\sigma_\gamma^2 \gamma_B^2 \sigma_\theta^2}{\sigma_\beta^4}$ , and  $\text{ARE}[\text{Var}(\hat{\gamma}_B)] = 1 - \frac{2\sigma_\gamma^2 \gamma_B^2 \sigma_\theta^2}{\sigma_\beta^4}$ . This demonstrates some gain in efficiency for both  $\gamma_X$  and  $\gamma_B$ .

### A.1.3 Approach 2: Constrained MLE

Depending on the information available from the external model  $Y|X$ , two possible situations correspond to two different constraints:

- *Approach 2.1: Only the estimated coefficient  $\beta$  is known from model 2.5*

From model 2.5–2.7, it is easy to see that the constraint is  $\theta = \frac{\beta - \gamma_X}{\gamma_B}$ , describing the relationship between the unknown variable  $\theta$ , the known variable  $\beta$  and the target variable  $\gamma$ . The log-likelihood is given by

$$\begin{aligned} l &= l(\gamma, \theta, \sigma_\gamma^2, \sigma_\theta^2) \\ &= \sum_{i=1}^n \log f(Y_i, B_i | X_i) = \sum_{i=1}^n \log f(Y_i | X_i, B_i; \gamma, \sigma_\gamma^2) + \sum_{i=1}^n \log f(B_i | X_i; \theta, \sigma_\theta^2) \\ &= -\frac{n}{2} \log(\sigma_\gamma^2) - \frac{1}{2\sigma_\gamma^2} \sum_{i=1}^n (Y_i - \gamma_X X_i - \gamma_B B_i)^2 - \frac{n}{2} \log(\sigma_\theta^2) - \frac{1}{2\sigma_\theta^2} \sum_{i=1}^n (B_i - \theta X_i)^2 \end{aligned} \quad (\text{A.4})$$

The goal is to maximize the log-likelihood A.4 over  $\gamma$ ,  $\sigma_\gamma$  and  $\sigma_\theta$  subject to the constraint  $\theta = \theta^* = \frac{\beta - \gamma_X}{\gamma_B}$ . By replacing  $\theta$  with  $\theta^*$ , taking the second derivative over  $\gamma$ , and taking the inverse of the information matrix, we obtain the asymptotic variance of  $\hat{\gamma}$ :

$$\text{Var}(\hat{\gamma}) = \mathbf{I}^{-1} = \frac{1}{n} \frac{\sigma_\gamma^2}{\sigma_\theta^2} \begin{pmatrix} \theta^{*2} + \frac{\sigma_\theta^4 \gamma_B^2}{\sigma_\gamma^2 + \sigma_\theta^2 \gamma_B^2} \sigma_X^{-2} & -\theta^* \\ -\theta^* & 1 \end{pmatrix}. \quad (\text{A.5})$$

Thus, the ARE of  $\text{Var}(\hat{\gamma})$  from the constrained MLE compared to the standard MLE is

$$\text{ARE}[\text{Var}(\hat{\gamma})] = \begin{pmatrix} 1 - \frac{\sigma_\theta^2}{\sigma_\theta^2 + \sigma_X^2} \frac{\sigma_\gamma^2}{\sigma_\gamma^2 + \gamma_B^2 \sigma_\theta^2} \\ 1 \end{pmatrix},$$

where we notice that there is some gain in efficiency for  $\gamma_X$  but no gain in efficiency for  $\gamma_B$ . We can see that the largest gain in efficiency is when  $\gamma_B$ ,  $\theta$  and  $\sigma_X$  are small.

- *Approach 2.2: Both of the estimated coefficient  $\beta$  and the standard deviation  $\sigma_\beta$  are known from model 2.5*

In this situation, knowing the true  $\sigma_\beta$  gives us more information which is incorporated through an additional constraint. In addition to the constraint  $\theta = \theta^* = \frac{\beta - \gamma_X}{\gamma_B^*}$  derived in approach 2.1, we add another constraint  $\text{Var}(Y|X) = \sigma_\beta^2 = \gamma_B^2 \sigma_\theta^2 + \sigma_\gamma^2$ , where  $\sigma_\theta^2 = \sigma_\theta^{*2} = \frac{\sigma_\beta^{*2} - \sigma_\gamma^2}{\gamma_B^2}$ . Then we maximize log-likelihood A.4 with respect to  $\gamma$  and  $\sigma_\gamma^2$  at fixed values  $\sigma_\theta^2 = \sigma_\theta^{*2}$ ,  $\theta = \theta^*$ . Note that different from approach 2.1,  $\sigma_\gamma^2$  and  $\gamma$  are not independent anymore. Thus, we need to consider  $\sigma_\gamma^2$  in the information matrix, and take the inverse of a  $3 \times 3$  matrix to get the correct asymptotic variance. Let  $\phi = (\gamma, \sigma_\gamma^2)^T$ ,

$$\mathbf{I} = -\text{E}_{\text{XB}} \left( \frac{\partial^2 l}{\partial \phi \phi^T} \right) = n \begin{pmatrix} \left( \frac{1}{\sigma_\gamma^2} + \frac{1}{\gamma_B^2 \sigma_\theta^{*2}} \right) \sigma_X^2 & \left( \frac{1}{\sigma_\gamma^2} + \frac{1}{\gamma_B^2 \sigma_\theta^{*2}} \right) \sigma_X^2 \theta^* & 0 \\ \left( \frac{1}{\sigma_\gamma^2} + \frac{1}{\gamma_B^2 \sigma_\theta^{*2}} \right) \sigma_X^2 \theta^* & \left( \frac{1}{\sigma_\gamma^2} + \frac{1}{\gamma_B^2 \sigma_\theta^{*2}} \right) (\sigma_\theta^{*2} + \sigma_X^2 \theta^{*2}) + \frac{1}{\gamma_B^2} & \frac{1}{\sigma_\theta^{*2} \gamma_B^3} \\ 0 & \frac{1}{\sigma_\theta^{*2} \gamma_B^3} & \frac{1}{2} \left( \frac{1}{\sigma_\gamma^4} + \frac{1}{\gamma_B^4 \sigma_\theta^{*4}} \right) \end{pmatrix},$$

By taking the inverse of  $\mathbf{I}$ , we can get the asymptotic variance of  $\hat{\gamma}$ :

$$\begin{cases} \text{Var}(\hat{\gamma}_B) &= \frac{1}{n} \frac{\sigma_\gamma^2}{\sigma_\theta^{*2}} \frac{\sigma_\gamma^4 + \gamma_B^4 \sigma_\theta^{*4}}{(\sigma_\gamma^2 + \gamma_B^2 \sigma_\theta^{*2})^2} = \frac{1}{n} \frac{\gamma_B^2 \sigma_\gamma^2}{\sigma_\beta^{*2} - \sigma_\gamma^2} \frac{\sigma_\gamma^4 + (\sigma_\beta^{*2} - \sigma_\gamma^2)^2}{\sigma_\beta^{*4}} \\ \text{Var}(\hat{\gamma}_X) &= \frac{1}{n} \frac{\sigma_\gamma^2}{\sigma_\theta^{*2}} \frac{1}{(\sigma_\gamma^2 + \gamma_B^2 \sigma_\theta^{*2})^2} \left[ (\sigma_\theta^{*2} \sigma_X^{-2} + \theta^{*2}) (\sigma_\gamma^4 + \gamma_B^4 \sigma_\theta^{*4}) - (\sigma_\gamma^2 - \gamma_B^2 \sigma_\theta^{*2}) \sigma_\gamma^2 \sigma_\theta^{*2} \sigma_X^{-2} \right] \\ &= (\sigma_\theta^{*2} \sigma_X^{-2} + \theta^{*2}) \text{Var}(\hat{\gamma}_B) - \frac{1}{n} \sigma_\gamma^4 \sigma_X^{-2} \frac{\sigma_\gamma^2 - \gamma_B^2 \sigma_\theta^{*2}}{\sigma_\beta^{*4}} \end{cases}$$

Thus, we can obtain the identical ARE to the synthetic data method (approach 1 in Appendix A.1.2). This demonstrates the asymptotic equivalence of the synthetic data approach with large  $m$  compared to the constrained ML approach that uses knowledge of all the parameters in the  $Y|X$  distribution.

### A.1.4 Approach 3: CSPML

This approach assumes that  $\beta$  is known, but does not assume that  $\sigma_\beta$  is known. To calculate the asymptotic variance of  $\hat{\gamma}$  in this approach, we need three matrices  $\mathbf{I}$ ,  $\mathbf{C}$  and  $\mathbf{L}$ . After some algebra,

it can be shown that  $\mathbf{C} = \frac{\sigma_X^2}{\sigma_\beta^2}(1, \theta^*)^T$ ,  $\mathbf{L} = \frac{n\gamma_B^2\sigma_\theta^*\sigma_X^2}{\sigma_\beta^4}$ . Thus,

$$\text{Cov}(\hat{\gamma}) = (\mathbf{I} + \mathbf{C}\mathbf{L}^{-1}\mathbf{C}^T)^{-1} = \frac{1}{n} \frac{\sigma_\gamma^2}{\sigma_\theta^2} \begin{pmatrix} \theta^{*2} + \frac{\sigma_\theta^4\gamma_B^2}{\sigma_\gamma^2 + \sigma_\theta^2\gamma_B^2}\sigma_X^{-2} & -\theta^* \\ -\theta^* & 1 \end{pmatrix},$$

which is identical to approach 2.1 in Appendix A.1.3.

## A.2 Derivation of asymptotic variances for special case 2: binary

### A.2.1 Standard MLE

We will be using the following notation:  $S \equiv S_\gamma(X, B) = \gamma_0 + \gamma_X X + \gamma_B B + \gamma_{XB} XB$ ,  $M \equiv M_\beta(X) = \beta_0 + \beta_1 X$  and  $K \equiv K_\theta(X) = \theta_0 + \theta_1 X$ . The ML estimators are the solution of maximizing  $\prod_{i=1}^n f(Y_i|X_i, B_i; \gamma)$ , i.e.  $\max_{\gamma} \left\{ \sum_{i=1}^n [Y_i S_i - \log\{1 + \exp(S_i)\}] \right\}$ .

In the tri-binary case, since  $X = X^2$ , and  $B = B^2$ , the Fisher information matrix has the following form:

$$\mathbf{I} = E_{XB} \mathbf{I}(X, B) = E_{XB} \left\{ \expit(S)[1 - \expit(S)] \begin{pmatrix} 1 & X & B & XB \\ X & X & XB & XB \\ B & XB & B & XB \\ XB & XB & XB & XB \end{pmatrix} \right\}.$$

There are four possible combinations of binary  $(X, B)$ , i.e.  $(0,0)$ ,  $(0,1)$ ,  $(1,0)$  and  $(1,1)$ . Thus, the expectation terms of the matrix  $\mathbf{I}(X, B)$  can be obtained through  $\frac{1}{n} \sum_{a,b \in \{0,1\}} \mathbf{I}(a, b) P(X = a, B = b)$ .

The asymptotic variance of  $\hat{\gamma}$  is given by:

$$\begin{cases} \text{Var}(\hat{\gamma}_0) = \frac{1}{n} \left( \frac{1}{P(BXY=001)} + \frac{1}{P(BXY=000)} \right), \\ \text{Var}(\hat{\gamma}_X) = \frac{1}{n} \sum_{a,b \in \{0,1\}} \frac{1}{P(BXY=0ab)}, \\ \text{Var}(\hat{\gamma}_B) = \frac{1}{n} \left( \frac{1}{P(B=1|XY=01)P(BXY=001)} + \frac{1}{P(B=1|XY=00)P(BXY=000)} \right), \\ \text{Var}(\hat{\gamma}_{XB}) = \frac{1}{n} \sum_{a,b \in \{0,1\}} \frac{1}{P(B=1|XY=ab)P(BXY=0ab)}, \end{cases}$$

where  $P(BXY = 0ab)$  is the probability of the  $(B=0, X=a, Y=b)$  combination, and  $P(B = 1|XY = ab)$  is the probability of  $B=1$  given  $X=a$  and  $Y=b$ ,  $a, b \in \{0, 1\}$ .



### A.2.2 Approach 1: Synthetic data method

Motivated by the MLE in the missing data problem (Little, 1992), we re-formate our target likelihood as follows:

$$\begin{aligned} \prod_{i=1}^n f(Y_i, X_i, B_i) \prod_{i=n+1}^{m+n} f(X_i, Y_i) &= \prod_{i=1}^n f(X_i, Y_i) f(B_i | X_i, Y_i) \prod_{i=n+1}^{m+n} f(X_i, Y_i) \\ &= \prod_{i=1}^n f(B_i | X_i, Y_i) \prod_{i=1}^{m+n} f(X_i, Y_i), \end{aligned} \quad (\text{A.6})$$

where  $f(B_i | X_i, Y_i)$  and  $f(X_i, Y_i)$  are independent from each other. The goal is to maximize likelihood A.6 over  $\gamma$ .

Let  $P(XY = ab) \equiv \Pr(X_i = a, Y_i = b)$ ,  $a, b \in \{0, 1\}$ ,  $i = 1, \dots, m+n$ . With the constraint  $\sum_{a,b \in \{0,1\}} P(XY = ab) = 1$ , there are three unknown variables in  $\prod_{i=1}^{m+n} f(X_i, Y_i)$ , i.e.  $P(XY=ab)$ ,  $a, b \in \{0, 1\}$ . Denote  $P(B = 1 | XY = ab) \equiv \Pr(B_i = 1 | X_i = a, Y_i = b)$ ,  $i = 1, \dots, n$ . Similarly, since there are four different combination of  $a$  and  $b$ , there are four unknown parameters in  $\prod_{i=1}^n f(B_i | X_i, Y_i)$ , i.e.  $P(B = 1 | XY = ab)$ , each of which is independent of others. By plugging in the four possible combinations of  $X$  and  $B$  into model 2.8 in Chapter 2, we can easily derive the expressions for  $\gamma$  as presented in Table A.1.

Table A.1: Formulas of  $\gamma$  in terms of  $P(B|XY)$  and  $P(XY)$

(X, B) combination	Transformation of model 2.8
(0, 0)	$\gamma_0 = \log \frac{P(B=0 XY=01)P(XY=01)}{P(B=0 XY=00)P(XY=00)}$
(1, 0)	$\gamma_0 + \gamma_X = \log \frac{P(B=0 XY=11)P(XY=11)}{P(B=0 XY=10)P(XY=10)}$
(0, 1)	$\gamma_0 + \gamma_B = \log \frac{P(B=1 XY=01)P(XY=01)}{P(B=1 XY=00)P(XY=00)}$
(1, 1)	$\gamma_0 + \gamma_X + \gamma_B + \gamma_{XB} = \log \frac{P(B=1 XY=11)P(XY=11)}{P(B=1 XY=10)P(XY=10)}$

Let  $mn(a, b)$  denote the number of observations with  $(X = a, Y = b)$  in the sample size of  $m+n$ . Since  $mn(a, b) \sim \text{Multinomial}(m+n, P(XY = ab))$ , we can easily obtain the MLE of  $P(XY)$ , and the corresponding estimated covariance as follows:

$$\begin{cases} \hat{P}(XY = ab) = \frac{mn(a,b)}{m+n} \\ \text{Var}[\hat{P}(XY = ab)] = \hat{P}(XY = ab)[1 - \hat{P}(XY = ab)] \\ \text{Cov}[\hat{P}(XY = ab), \hat{P}(XY = a'b')] = -\hat{P}(XY = ab)\hat{P}(XY = a'b') \end{cases}$$

Let  $n(a,b)$  be the number of observations with  $(X = a, Y = b)$  in the sample of size  $n$ , and  $n(B = 1|XY = ab)$  be the count of  $B = 1$  given  $(X = a, Y = b)$ . Since  $n(B = 1|XY = ab) \sim \text{Bin}(n(a,b), P(B = 1|XY = ab))$ , the MLE of  $P(B|X, Y)$  and its estimated covariance can be expressed as:

$$\begin{cases} \hat{P}(B = 1|XY = ab) = \frac{n(B=1|XY=ab)}{n(a,b)} \\ \text{Var}(\hat{P}(B = 1|XY = ab)) = \frac{\hat{P}(B=1|XY=ab)\hat{P}(B=0|XY=ab)}{n(a,b)} \\ \text{Cov}(\hat{P}(B = 1|XY = ab), \hat{P}(B = 1|XY = a'b')) = 0 \end{cases}$$

Therefore, the MLE of  $\gamma$  can be expressed as:

$$\begin{cases} \hat{\gamma}_0 = \log \frac{\hat{P}(B=0|XY=01) \hat{P}(XY=01)}{\hat{P}(B=0|XY=00) \hat{P}(XY=00)} \\ \hat{\gamma}_X = \log \frac{\hat{P}(B=0|XY=11) \hat{P}(B=0|XY=00) \hat{P}(XY=11) \hat{P}(XY=00)}{\hat{P}(B=0|XY=10) \hat{P}(B=0|XY=01) \hat{P}(XY=10) \hat{P}(XY=01)} \\ \hat{\gamma}_B = \log \frac{\hat{P}(B=1|XY=01) \hat{P}(B=0|XY=00)}{\hat{P}(B=0|XY=01) \hat{P}(B=1|XY=00)} \\ \hat{\gamma}_{XB} = \log \frac{\hat{P}(B=1|XY=11) \hat{P}(B=1|XY=00) \hat{P}(B=0|XY=10) \hat{P}(B=0|XY=01)}{\hat{P}(B=1|XY=10) \hat{P}(B=1|XY=01) \hat{P}(B=0|XY=00) \hat{P}(B=0|XY=11)} \end{cases}$$

According to the delta method and replacing the estimated proportions by the corresponding probabilities, we obtain the asymptotic variances as follow:

$$\begin{cases} \text{Var}(\hat{\gamma}_0) = \frac{1}{n} \sum_{(a,b) \in \{(0,1), (0,0)\}} \frac{P(B=1|XY=ab)}{P(BXY=0ab)} + \frac{1}{m+n} \sum_{(a,b) \in \{(0,1), (0,0)\}} \frac{1}{P(XY=ab)} \\ \text{Var}(\hat{\gamma}_X) = \frac{1}{n} \sum_{a,b \in \{0,1\}} \frac{P(B=1|XY=ab)}{P(BXY=0ab)} + \frac{1}{m+n} \sum_{a,b \in \{0,1\}} \frac{1}{P(XY=ab)} \\ \text{Var}(\hat{\gamma}_B) = \frac{1}{n} \sum_{(a,b) \in \{(0,1), (0,0)\}} \frac{1}{P(B=1|XY=ab)P(BXY=0ab)} \\ \text{Var}(\hat{\gamma}_{XB}) = \frac{1}{n} \sum_{a,b \in \{0,1\}} \frac{1}{P(B=1|XY=ab)P(BXY=0ab)} \end{cases}$$

where  $P(BXY = 0ab) = P(B = 0|XY = ab)P(XY = ab)$ .

Therefore, we find that the ARE of  $\text{Var}(\hat{\gamma})$  by adding  $m$  synthetic data observations compared to the original dataset of size  $n$  is

$$\text{ARE}[\text{Var}(\hat{\gamma})] = \mathbf{1} - (1 - \lambda) \begin{pmatrix} \frac{\sum_{(a,b) \in \{(0,1), (0,0)\}} 1/P(XY=ab)}{\sum_{(a,b) \in \{(0,1), (0,0)\}} 1/P(BXY=0ab)} \\ \frac{\sum_{a,b \in \{0,1\}} 1/P(XY=ab)}{\sum_{a,b \in \{0,1\}} 1/P(BXY=0ab)} \\ 0 \\ 0 \end{pmatrix} \quad (\text{A.7})$$

### A.2.3 Approach 2: Constrained MLE

When the summary-level information from model 2.9 is available in the form of coefficient estimates  $\beta$ , the constrained MLE is the solution of maximizing  $\prod_{i=1}^n f(Y_i|X_i, B_i; \gamma)f(B_i|X_i; \theta)$  subject

to the constraint that  $\Pr(Y = 1|X = x; \beta) = \sum_{b=0}^1 \Pr(Y = 1|X = x, B = b; \gamma)\Pr(B = b|X = x; \theta)$ .

The log-likelihood can be rewritten as

$$\max_{\gamma} \left\{ \sum_{i=1}^n [Y_i S_i - \log\{1 + \exp(S_i)\} + B_i K_i - \log\{1 + \exp(K_i)\}] \right\}$$

From the constraints, we can write  $\theta$  as a function of  $\gamma$  in the following way:

$$\begin{cases} \theta_0(\gamma) &= \text{logit} \left\{ \frac{\expit(\beta_0) - \expit(\gamma_0)}{\expit(\gamma_0 + \gamma_B) - \expit(\gamma_0)} \right\} \\ \theta_1(\gamma) &= \text{logit} \left\{ \frac{\expit(\beta_0 + \beta_1) - \expit(\gamma_0 + \gamma_X)}{\expit(\gamma_0 + \gamma_X + \gamma_B + \gamma_{XB}) - \expit(\gamma_0 + \gamma_X)} \right\} - \theta_0(\gamma) \end{cases}$$

Then  $K$  becomes  $K_{\gamma}(X) = \theta_0(\gamma) + \theta_1(\gamma)X$ . Denote  $\begin{cases} \sigma_{\beta_{0j}} &\equiv \frac{\partial}{\partial \gamma_j} \theta_0(\gamma) \\ \sigma_{\beta_{1j}} &\equiv \frac{\partial}{\partial \gamma_j} \theta_1(\gamma) \end{cases}$  where  $j = 0, 1, 2, 3$ .

The asymptotic variance of  $\hat{\gamma}$  can be derived through the  $4 \times 4$  matrix  $\frac{1}{n} \{E_{XB}[E_{Y|XB}(\mathbf{u}_{\gamma} \mathbf{u}_{\gamma}^T)]\}^{-1}$ , where

$$\begin{aligned} \mathbf{u}_{\gamma} &= \frac{\partial}{\partial \gamma} \log\{f(Y, B|X; \gamma, \theta(\gamma))\} \\ &= \begin{pmatrix} (\sigma_{\beta_{00}} + \sigma_{\beta_{10}}X)(B - \expit(K)) + Y - \expit(S) \\ (\sigma_{\beta_{01}} + \sigma_{\beta_{11}}X)(B - \expit(K)) + (Y - \expit(S))X \\ (\sigma_{\beta_{02}} + \sigma_{\beta_{12}}X)(B - \expit(K)) + (Y - \expit(S))B \\ (\sigma_{\beta_{03}} + \sigma_{\beta_{13}}X)(B - \expit(K)) + (Y - \expit(S))XB \end{pmatrix} \end{aligned}$$

Since all  $Y, X$  and  $B$  are binary variables, there are eight possible combinations of  $(Y, X, B)$ . Thus, the expectation term in the matrix  $E(\mathbf{u}_{\gamma} \mathbf{u}_{\gamma}^T)$  can be obtained through

$$\frac{1}{n} \sum_{a,b,c \in \{0,1\}} \mathbf{u}_{\gamma} \mathbf{u}_{\gamma}^T P(Y = a, X = b, B = c).$$

A variation of the above approach is when the external summary information comes in the form of the predicted probability, for any given  $X$ , i.e. we are simply provided with  $\bar{P}(X_i) = \hat{\Pr}(Y_i = 1|X_i)$ . In such case, it is easy to construct an estimation method that uses the estimated probability as a constraint. In addition, in the special case being considered here where  $Y$  and  $X$  are binary, it is easy to see that knowing  $\bar{P}(0)$  and  $\bar{P}(1)$  is equivalent to knowing  $\beta_0 = \text{logit}[\bar{P}(0)]$  and  $\beta_1 = \text{logit}[\bar{P}(1)] - \text{logit}[\bar{P}(0)]$ , so this also fits into the above framework to obtain the asymptotic

variance of  $\hat{\gamma}$ .

### A.2.4 Approach 3: CSPML

By implementing the specific distribution into the given formulas for  $\mathbf{I}$ ,  $\mathbf{C}$  and  $\mathbf{L}$ , we find that  $\mathbf{I}$  is the same as the information matrix in approach 1 in Appendix A.2.2, where  $\mathbf{C}_{4 \times 2}$  is the first two columns of matrix  $\mathbf{I}_{4 \times 4}$ , and

$$\mathbf{L} = E_{XB} \left\{ \expit(S)[1 - \expit(M)] - \expit(M)[1 - \expit(S)] \begin{pmatrix} 1 & X \\ X & X^2 \end{pmatrix} \right\},$$

where  $S \equiv S_{\gamma}(X, B) = \gamma_0 + \gamma_X X + \gamma_B B + \gamma_{XB} XB$  and  $M \equiv M_{\beta}(X) = \beta_0 + \beta_1 X$ . The calculation of  $\mathbf{L}$  is simple under the situation where  $X$  and  $B$  are both binary. Then  $\mathbf{I}$ ,  $\mathbf{C}$  and  $\mathbf{L}$  can be combined to give the variance of  $\hat{\gamma}$ .

Although we do not write out the formulas of  $\text{ARE}[\text{Var}(\hat{\gamma})]$  for approaches 2 and 3 in Appendix A.2.3 and A.2.4, respectively, they are numerically identical to the equation A.7 when  $\lambda = 0$ .

## APPENDIX B

### Appendix of Chapter 3

#### B.1 Summary of Features and Assumptions

Figure B.1 summarizes the key features and assumptions required in the internal-data-only regression, CSPML, EB, and the proposed approaches.

In the CSPML approach, the estimated  $Y|X$  model from the external data is assumed to be the best fitted model in the class, but that class of models may not include the true model. A common example is that when the true  $Y|X, B$  model belongs to the generalized linear model (GLM) family, it does not necessarily imply the  $Y|X$  model belongs to the same class of GLMs (e.g. when  $Y|X, B$  is linear model,  $Y|X$  could still be linear model if  $B|X$  is linear; but when  $Y|X, B$  is logistic model, it is common to fit a logistic model for  $Y|X$ , but the truth is not a logistic model as collapsibility does not hold for the logit link).

As showed in Figure B.1, the CSPML approach implicitly requires the joint distribution of  $(Y, X, B)$  to be identical in the internal and the external population, so it requires more than the mean exchangeability assumption, i.e. it requires more than the conditional distribution of  $Y|X, B$  to be the same. The CSPML will essentially always be biased if the  $Y|X, B$  distributions differ. As demonstrated in simulations, the CSPML can also be biased if the two populations only differ in the  $(X, B)$  distribution. The bias can be small or it can be large depending on the details of the situation. The EB estimator ameliorates differences between the population distributions by down-weighting the CSPML estimator if the data suggests it is biased.

Because the EB will down-weight the CSPML estimator if the lack of transportability leads to a poor CSPML estimate, it is robust to departures from the assumption of full transportability in specific external populations. The EB estimator is a shrinkage estimator that posits an additional stochastic framework for the underlying true parameter  $\gamma \sim N(\gamma_0, A)$ . The stochastic framework connects the internal estimator  $\hat{\gamma}_I$  and the CSPML estimator  $\hat{\gamma}_{CSPML}$ , and the difference  $\hat{\gamma}_I - \hat{\gamma}_{CSPML}$  is a measure of the distributional similarity of the joint distribution  $(Y, X, B)$  between the internal and the external population.

Approach*	Estimator	Target (Y X, B; $\gamma$ )		Additional Assumptions	Bias of $\hat{\gamma}$	Precision gains in estimating X coefficients
		Internal Study	External Study			
Internal Data Only	$\hat{\gamma}_I$	Correctly specified	-	-	Asymptotically unbiased	-
CSPML	$\hat{\gamma}_{CSPML}$	Correctly specified	Same as the internal study	(X, B) is the same in the internal and external studies  (Y X; $\hat{\beta}$ ) is the best fitted model in the class, but the class of models may not include the true model	Potential for bias when the joint distribution of (Y, X, B) is different in the two populations, which could arise from various forms of misspecification. E.g. B X, Y X,B, Y X or B.	The variance of X coefficients will usually decrease even when the joint distribution of (Y, X, B) is different in the two populations
EB	$\hat{\gamma}_{EB}$ , a shrinkage estimator that can be viewed as a weighted average of $\hat{\gamma}_I$ and $\hat{\gamma}_{CSPML}$	Correctly specified	-	$\gamma \sim N(\gamma_0, A)$ , relating $\hat{\gamma}_I$ and $\hat{\gamma}_{CSPML}$ : The difference $\hat{\gamma}_I - \hat{\gamma}_{CSPML}$ is used to measure the distributional similarity of the joint distribution (Y, X, B) between the internal and the external population.	Asymptotically unbiased  Small to moderate bias in finite samples.  Applies less weight to $\hat{\gamma}_{CSPML}$ if the lack of transportability leads to a poor $\hat{\gamma}_{CSPML}$ .	The smaller the difference between $\hat{\gamma}_I$ and $\hat{\gamma}_{CSPML}$ , the larger the efficiency gain. Efficiency gains diminish as the difference of $\hat{\gamma}_I$ and $\hat{\gamma}_{CSPML}$ becomes large.
Proposed	OCWE $\hat{\gamma}_{OCWE}$	Same as the EB approach			A weighted average of $\hat{\gamma}_{EB}$ 's, where the same weight is applied to each coefficient within $\hat{\gamma}_{EB}$	
	SC-Learner $\hat{\gamma}_{SC-Learner}$				A weighted average of $\hat{\gamma}_{EB}$ 's, where a different weight is attached to each coefficient.	

\*All approaches assume the parametric settings:  $g[E(Y|X, B)]$  is linear in (X, B); and  $g[E(Y|X)]$  is linear in X, where  $g(\cdot)$  is the known link function.  
**Abbreviations:** CSPML=Constrained Semi-Parametric Maximum Likelihood; EB=Empirical Bayes;  $\gamma_0$ =the underlying true parameter of  $\gamma$ ; A=some covariance matrix of  $\gamma$ ; E=optimal covariace weighted estimator; SC-Learner=selective coefficient learner.

Figure B.1: Summary of features and assumptions required in the direct regression, CSPML, EB and the proposed approaches.

## B.2 Derivation of the asymptotic Distribution

Let  $S_i(\gamma)$  denote the score function of  $Y|X, B$  distribution, and  $u, B, C$  and  $L$  defined in equation 3.2 from Section 3.2.4. The MLE of internal data solves  $\frac{1}{n} \sum_{i=1}^n S_i(\gamma_I) = 0$  to obtain  $\hat{\gamma}_I$ . CSPML solves  $\frac{1}{n} \sum_{i=1}^n \left( S_i(\gamma_{CSPML}) + \frac{\frac{\partial u(\gamma_{CSPML})}{\partial \gamma_{CSPML}}}{1 - \lambda^T u(\gamma_{CSPML})} \lambda \right) = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$  to obtain  $\hat{\gamma}_{CSPML}$  and  $\hat{\lambda}$ . We obtain the asymptotic distribution of target parameter  $\eta = (\gamma_{CSPML}, \lambda, \gamma_I)^T$  by the following steps:

1. Re-write the estimating equation  $g(\boldsymbol{\eta}) = \frac{1}{n} \sum_{i=1}^n \begin{pmatrix} S_i(\boldsymbol{\gamma}_{\text{CSPML}}) + \frac{\frac{\partial u(\boldsymbol{\gamma}_{\text{CSPML}})}{\partial \boldsymbol{\gamma}_{\text{CSPML}}}}{1 - \lambda^T u(\boldsymbol{\gamma}_{\text{CSPML}})} \lambda \\ \frac{u(\boldsymbol{\gamma}_{\text{CSPML}})}{1 - \lambda^T u(\boldsymbol{\gamma}_{\text{CSPML}})} \\ S_i(\boldsymbol{\gamma}_I) \end{pmatrix}$

2. Apply the first order Taylor series expansion to  $g(\boldsymbol{\eta})$  at true values  $(\boldsymbol{\gamma}_0, 0, \boldsymbol{\gamma}_0)^T$ :

$$\frac{1}{n} \sum_{i=1}^n \begin{pmatrix} S_i(\boldsymbol{\gamma}_0) \\ u_i(\boldsymbol{\gamma}_0) \\ S_i(\boldsymbol{\gamma}_0) \end{pmatrix} + \begin{pmatrix} -B & C & 0 \\ C & L & 0 \\ 0 & 0 & -B \end{pmatrix} \begin{pmatrix} \boldsymbol{\gamma}_{\text{CSPML}} - \boldsymbol{\gamma}_0 \\ \lambda - 0 \\ \boldsymbol{\gamma}_I - \boldsymbol{\gamma}_0 \end{pmatrix} + \text{op}(1) = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$$

$\Leftrightarrow$

$$\sqrt{n} \begin{pmatrix} \boldsymbol{\gamma}_{\text{CSPML}} - \boldsymbol{\gamma}_0 \\ \lambda - 0 \\ \boldsymbol{\gamma}_I - \boldsymbol{\gamma}_0 \end{pmatrix} = - \begin{pmatrix} -B & C & 0 \\ C & L & 0 \\ 0 & 0 & -B \end{pmatrix}^{-1} \frac{1}{\sqrt{n}} \begin{pmatrix} \sum_{i=1}^n S_i(\boldsymbol{\gamma}_0) \\ \sum_{i=1}^n u_i(\boldsymbol{\gamma}_0) \\ \sum_{i=1}^n S_i(\boldsymbol{\gamma}_0) \end{pmatrix} + \text{op}(1)$$

3. Asymptotic covariance under two conditions:

- When the uncertainty of  $\hat{\boldsymbol{\beta}}$  is ignorable (i.e.  $\hat{\boldsymbol{\beta}}$  can be used to approximate the true  $\boldsymbol{\beta}_0$ )

$$\begin{aligned} \text{Cov}(\hat{\boldsymbol{\eta}}) &= \frac{1}{n} \begin{pmatrix} -B & C & 0 \\ C & L & 0 \\ 0 & 0 & -B \end{pmatrix}^{-1} \frac{1}{n} n E \left\{ \begin{pmatrix} S_i(\boldsymbol{\gamma}_0) \\ u_i(\boldsymbol{\gamma}_0) \\ S_i(\boldsymbol{\gamma}_0) \end{pmatrix} \begin{pmatrix} S_i(\boldsymbol{\gamma}_0) \\ u_i(\boldsymbol{\gamma}_0) \\ S_i(\boldsymbol{\gamma}_0) \end{pmatrix}^T \right\} \begin{pmatrix} -B & C & 0 \\ C & L & 0 \\ 0 & 0 & -B \end{pmatrix}^{-1} \\ &= \frac{1}{n} \begin{pmatrix} -B & C & 0 \\ C & L & 0 \\ 0 & 0 & -B \end{pmatrix}^{-1} \begin{pmatrix} B & 0 & B \\ 0 & L & 0 \\ B & 0 & B \end{pmatrix} \begin{pmatrix} -B & C & 0 \\ C & L & 0 \\ 0 & 0 & -B \end{pmatrix}^{-1} \\ &= \frac{1}{n} \begin{pmatrix} (B + C^T L^{-1} C)^{-1} & 0 & (B + C^T L^{-1} C)^{-1} \\ 0 & B(B + C^T L^{-1} C)^{-1} & -(C + B^T C^{-1} L)^{-1} \\ 0 & -(C + B^T C^{-1} L)^{-1} & B^{-1} \end{pmatrix} \end{aligned}$$

- When the uncertainty of  $\hat{\boldsymbol{\beta}}$  is not ignorable

Assume that the external  $\boldsymbol{\beta}$  estimator  $\hat{\boldsymbol{\beta}}$  is derived from a finite sample of size  $N$ , whose variance is  $\frac{V_{\boldsymbol{\beta}}}{N}$ . We incorporate the uncertainty from  $\hat{\boldsymbol{\beta}}$  by applying the first order Taylor series expansion at the true  $\boldsymbol{\beta}_0$  to the only  $\boldsymbol{\beta}$ -related parameter  $u_i(\boldsymbol{\gamma}_0)$  in step 2. Let  $\rho = \lim_{N \rightarrow \infty} \frac{n}{N}$ , and  $Q = E\left\{ \frac{\partial u_i(\boldsymbol{\gamma}_0, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \right\}$

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n u_i(\boldsymbol{\gamma}_0, \hat{\boldsymbol{\beta}}) = \frac{1}{\sqrt{n}} \sum_{i=1}^n u_i(\boldsymbol{\gamma}_0, \boldsymbol{\beta}_0) + \sqrt{\rho} Q \sqrt{N} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)$$

Thus,

$$\begin{aligned}
\text{Cov}(\hat{\boldsymbol{\eta}}) &= \frac{1}{n} \begin{pmatrix} -B & C & 0 \\ C & L & 0 \\ 0 & 0 & -B \end{pmatrix}^{-1} \begin{pmatrix} B & 0 & B \\ 0 & L + \rho N Q^T \text{Var}(\hat{\boldsymbol{\beta}}) Q & 0 \\ B & 0 & B \end{pmatrix} \begin{pmatrix} -B & C & 0 \\ C & L & 0 \\ 0 & 0 & -B \end{pmatrix}^{-1} \\
&= \frac{1}{n} \begin{pmatrix} -B & C & 0 \\ C & L & 0 \\ 0 & 0 & -B \end{pmatrix}^{-1} \begin{pmatrix} B & 0 & B \\ 0 & L + \rho Q^T V_{\beta} Q & 0 \\ B & 0 & B \end{pmatrix} \begin{pmatrix} -B & C & 0 \\ C & L & 0 \\ 0 & 0 & -B \end{pmatrix}^{-1} \\
&= \frac{1}{n} \begin{pmatrix} \text{Var}(\hat{\boldsymbol{\gamma}}_{\text{CSPML}})[1 + \rho(C^T L^{-1} Q^T V_{\beta} Q L^{-1} C) \text{Var}(\hat{\boldsymbol{\gamma}}_{\text{CSPML}})] & \text{term}_1 & \text{Var}(\hat{\boldsymbol{\gamma}}_{\text{CSPML}}) \\ & \text{term}_2 & -(C + B^T C^{-1} L)^{-1} \\ & & B^{-1} \end{pmatrix}
\end{aligned}$$

$$\text{where } \begin{cases} \text{Var}(\hat{\boldsymbol{\gamma}}_{\text{CSPML}}) = (B + C^T L^{-1} C)^{-1}, \\ \text{term}_1 = \rho \text{Var}(\hat{\boldsymbol{\gamma}}_{\text{CSPML}}) (B L^{-1} Q^T V_{\beta} Q L^{-1} C) \text{Var}(\hat{\boldsymbol{\gamma}}_{\text{CSPML}}), \\ \text{term}_2 = B \text{Var}(\hat{\boldsymbol{\gamma}}_{\text{CSPML}}) + B \text{Var}(\hat{\boldsymbol{\gamma}}_{\text{CSPML}}) (\rho Q^T V_{\beta} Q) \text{Var}(\hat{\boldsymbol{\gamma}}_{\text{CSPML}}) B. \end{cases}$$

### B.3 Re-parameterize $\hat{\boldsymbol{\gamma}}_{\text{EB}}$

Based on matrix 3.2 in *Proposition 1* of Section 3.2.4, we can show that

$$\begin{cases} \hat{\boldsymbol{\gamma}}_{\text{I}} \mid \boldsymbol{\gamma} \sim N(\boldsymbol{\gamma}, \hat{\Sigma}) \\ \hat{\boldsymbol{\gamma}}_{\text{CSPML}} \mid \boldsymbol{\gamma} \sim N(\boldsymbol{\gamma}, \hat{V}_{\text{CSPML}}) \end{cases} \implies \begin{pmatrix} Z \\ \hat{\boldsymbol{\gamma}}_{\text{CSPML}} \\ \hat{\boldsymbol{\gamma}}_{\text{I}} \end{pmatrix} \sim N \left( \begin{pmatrix} 0 \\ \boldsymbol{\gamma} \\ \boldsymbol{\gamma} \end{pmatrix}, \begin{pmatrix} \hat{V}_{\text{Z}} \equiv \hat{V}_{\text{I}} - \hat{V}_{\text{CSPML}} & 0 & \hat{V}_{\text{Z}} \\ 0 & \hat{V}_{\text{CSPML}} & \hat{V}_{\text{CSPML}} \\ \hat{V}_{\text{Z}} & \hat{V}_{\text{CSPML}} & \hat{V}_{\text{I}} \end{pmatrix} \right).$$

$\text{Cov}(Z, \hat{\boldsymbol{\gamma}}_{\text{CSPML}}) = 0$  can also be derived from the conclusion that  $\text{Cov}(\hat{\boldsymbol{\gamma}}_{\text{I}}, \hat{\boldsymbol{\gamma}}_{\text{CSPML}}) = \text{Var}(\hat{\boldsymbol{\gamma}}_{\text{CSPML}})$  in *Proposition 1*. Thus, we can use  $Z$ , the scalar  $Z^T V_{\text{I}}^{-1} Z$  and Woodbury formula [Golub and Loan, 1996] to re-parameterize  $\hat{\boldsymbol{\gamma}}_{\text{EB}}$ :

$$\begin{aligned}
\hat{\boldsymbol{\gamma}}_{\text{EB}} &= \hat{A}(\hat{\Sigma} + \hat{A})^{-1} \hat{\boldsymbol{\gamma}}_{\text{I}} + \hat{\Sigma}(\hat{\Sigma} + \hat{A})^{-1} \hat{\boldsymbol{\gamma}}_{\text{CSPML}} \\
&= Z Z^T (\hat{V}_{\text{I}} + Z Z^T)^{-1} (Z + \hat{\boldsymbol{\gamma}}_{\text{CSPML}}) + \hat{V}_{\text{I}} (\hat{V}_{\text{I}} + Z Z^T)^{-1} \hat{\boldsymbol{\gamma}}_{\text{CSPML}} \\
&= \hat{\boldsymbol{\gamma}}_{\text{CSPML}} + Z Z^T (\hat{V}_{\text{I}} + Z Z^T)^{-1} Z \\
&= \hat{\boldsymbol{\gamma}}_{\text{CSPML}} + Z \left( 1 - \frac{1}{1 + Z^T V_{\text{I}}^{-1} Z} \right).
\end{aligned}$$

Or equivalently,  $\hat{\boldsymbol{\gamma}}_{\text{EB}} = \hat{\boldsymbol{\gamma}}_{\text{I}} - \hat{V}_{\text{I}} (\hat{V}_{\text{I}} + Z Z^T)^{-1} Z = \hat{\boldsymbol{\gamma}}_{\text{I}} - Z \frac{1}{1 + Z^T V_{\text{I}}^{-1} Z}$ .



## B.4 Additional Simulation Results

Figure B.2 summarizes the additional simulation settings we assessed regarding simulation III in Section 3.3. The goal is to evaluate the performance of proposed approach and the comparisons when the joint distribution of  $(X, B)$  is misspecified in the external model 1. Scenario (1) is the same as simulation III in the main manuscript; scenario (2) assesses different  $(X, B)$  in external model 1 through the different conditional distribution  $B|X$ ; and scenario (3) varies the joint distribution of  $(X, B)$  through the marginal distribution of  $X$ .

	Internal Data	External Model		
		1	2	3
Model	$Y X1, X2, X3, X4, B$	$Y X1, X2$	$Y X1, X2$	$Y X1, X2, X3, X4$
The joint distribution of $(X, B)$	$(X, B) \sim N((0,0), 1)$ , correlation 0.3	(1) $(X, B)$ is different from the marginal distribution of $B$ : $(X, B) \sim N((0, 1.5), 1)$ , correlation 0.3		
	$X \sim N(0, 1)$ , correlation 0.3;	(2) $(X, B)$ is different from the conditional distribution $B X$ : $X \sim N(0, 1)$ , correlation 0.3; $B X \sim (-1+0.25X, 1)$	Same as the internal data	Same as the internal data
	$B X \sim (-1+0.5X, 1)$	(3) $(X, B)$ is different from the marginal distribution $X$ : $X \sim N(0.25, 1)$ , correlation 0.3; $B X \sim (-1+0.5X, 1)$		
Sample size	$n=200$	$m_1=30,000$	$m_2=30,000$	$m_3=30,000$

Figure B.2: Simulation settings of additional scenarios for simulation III in Chapter 3: the joint distribution of  $(X, B)$  is misspecified in the external model 1.

For simulation IV in Section 3.3, we have additionally assessed different forms of misspecified outcome model in the external population listed in Figure B.3, including different intercept in scenario (1), different  $X$  coefficients in scenario (2), and different intercept and  $X$  coefficients in scenario (3). Scenario (3) is the same as simulation IV in the main manuscript.

Internal Data		External Model		
		1	2	3
Model	$Y X_1, X_2, X_3, X_4, B$	$Y X_1, X_2$	$Y X_1, X_3$	$Y X_1, X_2, X_3, X_4$
Outcome model	$\text{logit Pr}(Y=1 X, B) = -1 - 0.5X + 0.5B$	Same as the internal data	Same as the internal data	(1) intercept different: $\text{logit Pr}(Y=1 X, B) = 1 - 0.5X + 0.5B$
				(2) X coefficients different: $\text{logit Pr}(Y=1 X, B) = -1 + 0.5X + 0.5B$
				(3) intercept & X coefficients different: $\text{logit Pr}(Y=1 X, B) = 1 + 0.5X + 0.5B$
Sample size	$n=200$	$m_1=30,000$	$m_2=30,000$	$m_3=30,000$
<b>Note:</b> (X, B) followed Gaussian distribution with mean zero, standard deviation 1, and correlation 0.3.				

Figure B.3: Simulation settings of additional scenarios for simulation IV in Chapter 3: the outcome model is misspecified in the external model 3.

The results in Figure B.4 corresponding to Figure B.2 indicate that CSPML estimators would have substantial bias when the difference comes from the conditional distribution  $B|X$  or marginal distribution  $B$ , but not troublesome when the marginal distribution  $X$  is different. In the cases where CSPML estimators are biased, the EB estimator will protect against the bias by downweighting the CSPML estimator and sacrificing some precision. In summary, these simulation results provided the evidence that the proposed approach is robust to the heterogeneous covariate distribution (X, B) of the external population.

The result in Figure B.5 corresponding to Figure B.3 indicates that when the misspecification gets worse (from misspecified intercept only to both misspecified intercept and X coefficients), the EB estimator will trade more efficiency gain to correct the bias from the CSPML estimator and thus we obtain less efficiency by incorporating this external model (OCWE assigned smaller weights on the corresponding estimator  $EB_3$ ). More specifically, when we had mild misspecification in the external model 3 (only the intercept, i.e., the prevalence in the external population is different from the internal population), the external model 3 can still contribute 40% in the final OCWE. However, when the misspecification gets worse (both intercept and X coefficients are misspecified), the OCWE basically rejects the wrong information from the external model 3 by incorporating only 4% of it.

## B.5 External Calculators in Prostate Cancer Data Example

In the real data analysis in Section 3.4, we incorporate two external risk calculators of high-grade prostate cancer, the Prostate Cancer Prevention Trial (PCPT<sub>hg</sub>) [Thompson et al., 2006] and the European Randomized Study of Screening for Prostate Cancer (ERSPC) risk calculator 3 [Roobol

	Internal Data + External 1		Internal Data + External 2		Internal Data + External 3			Composite of EB Estimators		SC-Learner
	CSPML 1		CSPML 2		CSPML 3			IVW	OCWE	
	EB 1	EB 2	EB 3	EB 3	EB 3	EB 3	EB 3			
Direct Regression										
Bias (SD) [ESE]										
95% Coverage Rate										
(1)	Weights	/	/	/	/	/	/	[.33, .33, .33]	[.12, .38, .49]	/
	Y <sub>0</sub>	.665 (.080) [.102] 0%	-0.01 (.201) [.162] 90%	-0.39 (.081) [.085] 97%	-0.40 (.163) [.161] 95%	-0.33 (.051) [.058] 0%	-0.39 (.169) [.169] 93%	-0.27 (.177) [.163] 93%	-0.31 (.175) [.163] 93%	-0.27 (.177) [.163] 93%
	Y <sub>1</sub>	.02 (.094) [.112] 98%	-0.18 (.219) [.175] 89%	.012 (.095) [.097] 97%	-0.11 (.190) [.174] 93%	.0220 (.066) [.074] 100%	-0.13 (.198) [.182] 90%	-0.14 (.202) [.176] 92%	-0.12 (.197) [.177] 93%	-0.14 (.202) [.176] 92%
	Y <sub>2</sub>	.004 (.090) [.112] 98%	-0.16 (.210) [.175] 90%	-0.15 (.091) [.096] 97%	-0.15 (.180) [.173] 94%	.005 (.062) [.073] 100%	-0.15 (.188) [.182] 91%	-0.16 (.192) [.176] 93%	-0.15 (.187) [.177] 94%	-0.16 (.199) [.177] 93%
	Y <sub>3</sub>					-0.14 (.062) [.072] 100%	-0.17 (.178) [.194] 93%	-0.18 (.198) [.195] 96%	-0.17 (.193) [.192] 96%	-0.18 (.192) [.192] 96%
	Y <sub>4</sub>					-0.21 (.067) [.073] 97%	-0.25 (.188) [.195] 94%	-0.26 (.207) [.196] 94%	-0.25 (.202) [.193] 94%	-0.25 (.188) [.195] 96%
	Y <sub>B</sub>									
	(2)	Weights	/	/	/	/	/	/	[.33, .33, .33]	[.26, .26, .47]
Y <sub>0</sub>		.055 (.167) [.164] 94%	-0.09 (.243) [.222] 94%	-0.25 (.167) [.163] 95%	-0.23 (.235) [.222] 96%	-0.13 (.160) [.154] 95%	-0.21 (.244) [.229] 93%	-0.17 (.240) [.224] 95%	-0.17 (.240) [.225] 95%	-0.17 (.240) [.224] 95%
Y <sub>1</sub>		-.169 (.140) [.143] 83%	-0.07 (.221) [.208] 93%	-0.10 (.139) [.142] 95%	-0.38 (.212) [.208] 94%	.007 (.122) [.127] 95%	-0.38 (.217) [.216] 94%	-0.49 (.210) [.210] 94%	-0.47 (.211) [.211] 94%	-0.49 (.210) [.210] 94%
Y <sub>2</sub>		-.174 (.147) [.144] 82%	-0.65 (.231) [.209] 93%	-0.41 (.146) [.142] 95%	-0.41 (.223) [.209] 93%	-0.29 (.127) [.127] 96%	-0.39 (.228) [.216] 94%	-0.48 (.210) [.210] 94%	-0.47 (.211) [.211] 94%	-0.52 (.211) [.211] 94%
Y <sub>3</sub>						-0.17 (.126) [.126] 96%	-0.22 (.227) [.228] 95%	-0.22 (.210) [.229] 93%	-0.22 (.211) [.226] 93%	-0.22 (.211) [.226] 93%
Y <sub>4</sub>						-0.55 (.124) [.126] 96%	-0.21 (.206) [.227] 97%	-0.16 (.223) [.228] 96%	-0.17 (.219) [.225] 96%	-0.21 (.206) [.227] 97%
Y <sub>B</sub>										
(3)		Weights	/	/	/	/	/	/	[.33, .33, .33]	[.26, .27, .47]
	Y <sub>0</sub>	-.061 (.167) [.163] 94%	-0.30 (.235) [.222] 95%	-0.25 (.167) [.163] 95%	-0.23 (.235) [.222] 96%	-0.13 (.160) [.154] 95%	-0.21 (.244) [.229] 93%	-0.25 (.238) [.224] 95%	-0.24 (.239) [.225] 95%	-0.25 (.238) [.224] 95%
	Y <sub>1</sub>	-.008 (.139) [.142] 95%	-0.37 (.212) [.208] 94%	-0.10 (.139) [.142] 95%	-0.38 (.212) [.208] 94%	.007 (.122) [.127] 95%	-0.38 (.217) [.216] 94%	-0.38 (.213) [.210] 94%	-0.38 (.214) [.211] 94%	-0.38 (.213) [.210] 94%
	Y <sub>2</sub>	-.030 (.146) [.142] 95%	-0.38 (.224) [.240] 93%	-0.41 (.146) [.142] 95%	-0.41 (.223) [.209] 93%	-0.29 (.127) [.127] 96%	-0.39 (.228) [.216] 94%	-0.39 (.224) [.210] 94%	-0.39 (.225) [.211] 94%	-0.38 (.225) [.211] 94%
	Y <sub>3</sub>					-0.17 (.126) [.126] 96%	-0.22 (.227) [.228] 95%	-0.22 (.245) [.229] 93%	-0.22 (.241) [.226] 93%	-0.22 (.240) [.226] 93%
	Y <sub>4</sub>					-0.55 (.124) [.126] 96%	-0.21 (.206) [.227] 97%	-0.16 (.223) [.228] 96%	-0.16 (.220) [.225] 95%	-0.21 (.206) [.227] 97%
	Y <sub>B</sub>									
	<b>Abbreviations:</b> CSPML, constrained maximal likelihood estimator; EB, empirical Bayes estimator; SD, Monte Carlo standard deviation from 500 simulations; ESE, estimated average standard error from 500 simulations; IVW, Inverse variance-weighted estimator; OCWE, Optimal covariance-weighted estimator; SC-Learner, Selective coefficient learner. <b>Note:</b> Internal dataset size n=200; green represents good performance with small bias, yellow represents overestimated ESE (in source brackets) compared with the Monte Carlo SD (in round brackets), and red represents the poor performance of bias/95% coverage rate.									

**Abbreviations:** CSPML, constrained maximal likelihood estimator; EB, empirical Bayes estimator; SD, Monte Carlo standard deviation from 500 simulations; ESE, estimated average standard error from 500 simulations; IVW, Inverse variance-weighted estimator; OCWE, Optimal covariance-weighted estimator; SC-Learner, Selective coefficient learner. **Note:** Internal dataset size n=200; green represents good performance with small bias, yellow represents overestimated ESE (in square brackets) compared with the Monte Carlo SD (in round brackets), and red represents the poor performance of bias/95% coverage rate.

Figure B.4: Simulation results of additional scenarios of simulation III in Chapter 3.

Direct		Internal Data + External 1		Internal Data + External 2		Composite of EB Estimators				
Regression		CSPML 1	EB 1	CSPML 2	EB 2	CSPML 3	EB 3	IVW	OCWE	SC-Learner
Bias (SD) [ESE]										
95% Coverage Rate										
Y <sub>0</sub>	-041 (198) [193]	-008 (.081) [085]	-032 (.163) [161]	-007 (.089) [084]	-032 (.164) [161]	1.89 (.054) [142]	-026 (.198) [174]	[33, 33, 33]	[-41, 17, 41]	/
	95%	96%	94%	95%	95%	0%	93%	-034 (.174) [164]	-029 (.177) [165]	-035 (.173) [163]
	-020 (.230) [208]	-005 (.095) [097]	-014 (.190) [174]	-002 (.089) [097]	-014 (.189) [174]	.006 (.065) [149]	-019 (.227) [186]	-020 (.201) [177]	-018 (.203) [177]	-020 (.200) [176]
	94%	96%	93%	97%	93%	100%	90%	90%	90%	90%
Y <sub>1</sub>	-018 (.220) [208]	-008 (.091) [096]	-014 (.180) [175]			.010 (.064) [149]	-017 (.218) [187]	-017 (.205) [187]	-014 (.201) [182]	-017 (.197) [175]
	95%	97%	94%			100%	91%	91%	91%	91%
Y <sub>2</sub>	-019 (.208) [207]		-013 (.171) [180]	-003 (.089) [095]	-013 (.171) [180]	-012 (.064) [147]	-019 (.207) [185]	-019 (.195) [186]	-016 (.200) [190]	-018 (.188) [178]
	96%		96%	98%	96%	100%	93%	93%	93%	93%
Y <sub>3</sub>	-026 (.216) [208]					-012 (.066) [148]	-026 (.214) [189]	-026 (.216) [200]	-026 (.215) [198]	-026 (.214) [189]
	94%					97%	94%	94%	94%	94%
Y <sub>4</sub>	-019 (.218) [207]									
	98%									
Y <sub>B</sub>										
Y <sub>0</sub>						/	/	[34, 33, 32]	[54, 33, 13]	/
						-03 (.089) [153]	-04 (.196) [175]	-038 (.173) [164]	-035 (.165) [161]	-038 (.172) [163]
						99%	93%	94%	94%	94%
Y <sub>1</sub>						858 (.137) [164]	-013 (.227) [189]	-017 (.200) [177]	-016 (.193) [174]	-017 (.199) [177]
						0%	90%	92%	93%	92%
Y <sub>2</sub>						.864 (.136) [165]	-009 (.216) [189]	-013 (.203) [188]	-012 (.195) [186]	-012 (.195) [176]
						0%	92%	93%	94%	93%
Y <sub>3</sub>						855 (.138) [163]	-013 (.208) [188]	-016 (.195) [187]	-013 (.196) [190]	-016 (.188) [180]
						1%	94%	95%	95%	95%
Y <sub>4</sub>						875 (.139) [164]	-021 (.215) [191]	-025 (.216) [200]	-026 (.216) [205]	-021 (.215) [191]
						1%	93%	94%	94%	93%
Y <sub>B</sub>										
Y <sub>0</sub>						/	/	[34, 33, 32]	[57, 38, 04]	/
						1.69 (.189) [162]	-031 (.169) [178]	-032 (.173) [165]	-030 (.164) [160]	-032 (.172) [164]
						0%	96%	93%	93%	93%
						.747 (.178) [177]	-016 (.197) [192]	-015 (.201) [178]	-013 (.191) [173]	-015 (.200) [178]
Y <sub>1</sub>						3%	95%	92%	93%	92%
						.756 (.174) [177]	-014 (.188) [192]	-015 (.206) [189]	-013 (.196) [187]	-014 (.197) [177]
Y <sub>2</sub>						3%	94%	93%	94%	93%
						.762 (.167) [175]	-014 (.178) [191]	-015 (.194) [187]	-013 (.193) [189]	-014 (.187) [181]
Y <sub>3</sub>						1%	95%	95%	94%	95%
						.746 (.179) [176]	-022 (.188) [193]	-025 (.216) [202]	-026 (.217) [207]	-022 (.216) [193]
Y <sub>4</sub>						4%	93%	94%	95%	94%

**Abbreviations:** CSPML, constrained maximal likelihood estimator; EB, empirical Bayes estimator; SD, Monte Carlo standard deviation from 500 simulations; ESE, estimated average standard error from 500 simulations; IVW, Inverse variance-weighted estimator; OCWE, Optimal covariance-weighted estimator; SC-Learner, Selective coefficient learner. **Note:** Internal dataset size n=200; green represents good performance with small bias, yellow represents overestimated ESE (in square brackets) compared with the Monte Carlo SD (in round brackets), and red represents the poor performance of bias/95% coverage rate.

Figure B.5: Simulation results of additional scenarios of simulation IV in Chapter 3.

et al., 2012].

The proposed framework requires the consistent type of the common parameters shared across different external models (e.g. if one external model used mean-centered continuous age as a predictor, then all other external models including the target outcome model need to use the same age variable). Since PCPThg used log-based PSA and ERSPC used mean-centered log2-based PSA, we need to reconcile the different transformation by adjusting the reported intercept. In PCPThg, the authors only reported median PSA [Thompson et al., 2006] while ERSPC reported both median and mean PSA in their Table 1 [Roobol et al., 2012]. Therefore, we decided to use median-centered log2-based PSA, mean-centered continuous age, and mean-centered biopsy variable throughout.

We adjusted the originally reported estimated coefficients and intercept as follows:

- **Original PCPThg:**  $\text{logit}(p_i) = -6.25 + 1.29\log(\text{PSA}_i) + \text{DRE}_i + 0.03\text{Age}_i - 0.36\text{Biopsy}_i + 0.96\text{Race}_i$ ;
- $\hat{\beta}_{\log_2(\text{PSA})} = \frac{\hat{\beta}_{\log(\text{PSA})}}{\log_2(e)} = \frac{1.29}{\log_2(e)} = 0.8941599$ ;
- **Adjusted**  $\hat{\beta}_0 = \hat{\beta}_0 + \hat{\beta}_{\log_2(\text{PSA})} \times [\text{median } \log_2(\text{PSA})] + \hat{\beta}_{\text{Age}} \overline{\text{Age}} + \hat{\beta}_{\text{Biopsy}} \overline{\text{Biopsy}}$   
 $= -6.25 + 0.8941599 \times \log_2(1.5) + 0.03 \times 69.66 - 0.36 \times (753/5519)$   
 $= -3.686268$

where median PSA can be found in the first paragraph of Results, and mean age (i.e.  $\overline{\text{Age}}$ ) and mean biopsy (i.e.  $\overline{\text{Biopsy}}$ ) can be estimated from Table 2 in Thompson et al. [2006].

Similarly, since the  $\log_2(\text{PSA})$  in the original ERSPC risk calculator was mean-centered, we need to transform it to median-centered by adjusting the intercept as follows:

- **Original ERSPC:**  $\text{logit}(p_i) = -3.51 + 1.18\log_2(\text{PSA}_i) + 1.81\text{DRE}_i - 1.51\log_2(\text{TRUS} - \text{PV}_i)$
- **Adjusted**  $\hat{\beta}_0 = \hat{\beta}_0 + \hat{\beta}_{\log_2(\text{PSA})} \times [\text{median } \log_2(\text{PSA}) - \overline{\log(\text{PSA})}]$   
 $\approx \hat{\beta}_0 + \hat{\beta}_{\log_2(\text{PSA})} \times [\log_2(\text{median PSA}) - \log(\overline{\text{PSA}})]$   
 $= -3.51 + 1.18 \times [\log_2(4.3) - \log(6.1)]$   
 $= -3.16$

where mean and median PSA can be found in Table 1 from Roobol et al. [2012].

For other estimated coefficients that are not mentioned here, we used the originally reported values in the analysis.

## APPENDIX C

### Appendix of Chapter 4

#### C.1 Deriving the Initial Estimates for the External Populations

In this section, we will show how to obtain the initial parameter estimates of external population  $k$ . Let  $(\hat{\gamma}_0^{S_0}, \hat{\gamma}_X^{S_0T}, \hat{\gamma}_B^{S_0T})^T$  be the direct regression estimates of  $Y|X, B, S = 0$  using internal data only. For external population  $k$ , we know the parameter estimates  $\hat{\beta}_k = (\hat{\beta}_0, \hat{\beta}_X^T)^T$  from the fitted model  $Y|X_k; \beta_k$ . We assume that all predictors,  $X$  and  $B$ , are centered, and the true target model parameter for the external population  $k$  is  $(\gamma_0^{S_k}, \gamma_X^{S_kT}, \gamma_B^{S_0T})^T$ , assuming the coefficient of the unobserved variable  $B$  is the same as the internal population, i.e.  $\gamma_B^{S_k} = \gamma_B^{S_0}$ .

The goal of estimating  $\gamma_0^{S_k}$  and  $\gamma_X^{S_k}$  from model  $Y|X, B, S = k$ ;  $\gamma^{S_k}$  is equivalent to correcting the bias of  $\hat{\beta}_k$  in the reduced model  $Y|X_k; \beta_k$  considering covariates  $X_{(-k)}$  and  $B$  as omitted. To simplify notation, we assume  $B$  is the only omitted covariate in the derivation below. Neuhaus and Jewell [1993] provided a Taylor-series-expansion approximation to show that the ratio of coefficients remains constant in both the reduced and the full model when the omitted  $B$  is independent of the observed  $X$ , i.e.  $\frac{\gamma_{X_1}}{\gamma_{X_2}} \approx \frac{\beta_{X_1}}{\beta_{X_2}}$ , indicating that the relative effect size among regression coefficients remains consistent across models. In their Table 3 and equation 9, Neuhaus and Jewell [1993] provided the algebraic relationship between  $\gamma_X^{S_k}$  and  $\beta_X$  for exponential family when the omitted  $B$  and the observed  $X$  are correlated. In the subsequent paragraphs, we will explain in detail how to estimate  $\gamma_0^{S_k}$  and  $\gamma_X^{S_k}$  in linear regression (continuous  $Y$ ) and logistic regression (binary  $Y$ ), respectively.

**1. Linear Regression:** Suppose  $E(B|X; \theta) = \theta X$ . We start by replacing  $B$  with the conditional expected value  $E(B|X; \theta)$  in the mean profile of the target model:

$$E(Y|X, B; \gamma) = \gamma_0^{S_k} + \gamma_X^{S_kT} X + \gamma_B^{S_0T} B = \gamma_0^{S_k} + \gamma_X^{S_kT} X + \gamma_B^{S_0T} \theta X = E(Y|X; \gamma, \theta)$$

Since  $\hat{E}(Y|X; \beta) = \hat{\beta}_0 + \hat{\beta}_X^T X$  is available through the externally fitted model, we can obtain the estimation of  $\gamma_0^{S_k}$  and  $\gamma_X^{S_k}$  by matching the intercept and  $X$  coefficient between  $\hat{E}(Y|X; \gamma, \theta)$  and  $\hat{E}(Y|X; \beta)$ , respectively:  $\hat{\gamma}_0^{S_k} = \hat{\beta}_0$  and  $\hat{\gamma}_X^{S_k} = \hat{\beta}_X - \theta^T \hat{\gamma}_B^{S_0}$ . In a special case where the internal and



the external population only differ in intercept and  $S$  is independent of  $\mathbf{X}$  and  $\mathbf{B}$ , we can directly obtain the initial estimates  $\hat{\gamma}_k = (\hat{\beta}_0, \hat{\gamma}_X^{S_k^T}, \hat{\gamma}_B^{S_0^T})^T$ .

**2. Logistic Regression:** In logistic regression where  $g()$  is the logit link function, we connect the intercepts  $\beta_0$  and  $\gamma_0^{S_k}$  through the equation  $\text{logit}^{-1}(\beta_0) = E_{B|X}(\mu_0^{S_k})$ , where  $\mu_0^{S_k} = g^{-1}(Y|\mathbf{X}, \mathbf{B}; \gamma_X^{S_k} = 0) = \text{logit}^{-1}(\gamma_0^{S_k} + \mathbf{B}^T \gamma_B^{S_0})$ . For the right hand side, we expand  $\mathbf{B}$ , a vector of length  $Q$ , at  $E(\mathbf{B}|\mathbf{X})$  using the third-order Taylor series expansion as follows:

$$\begin{aligned} E_{B|X}(\mu_0^{S_k}) &= E_{B|X}[\text{logit}^{-1}(\gamma_0^{S_k} + \mathbf{B}^T \gamma_B^{S_0})] \\ &\approx \text{logit}^{-1}(w) \left\{ 1 + \frac{1}{2} \frac{1 - e^w}{(1 + e^w)^2} \sum_{i=1}^Q \sum_{j=1}^Q \gamma_{B_i}^{S_0} \gamma_{B_j}^{S_0} E_{B|X}[(B_i - E(B_i|\mathbf{X})) (B_j - E(B_j|\mathbf{X}))] \right\} \\ &= \text{logit}^{-1}(w) \left[ 1 + \frac{1}{2} \frac{1 - e^w}{(1 + e^w)^2} \text{Var} \left( \sum_{i=1}^Q \gamma_{B_i}^{S_0} B_i | \mathbf{X} \right) \right] \end{aligned} \quad (\text{C.1})$$

where  $w = \gamma_0^{S_k} + E(\mathbf{B}^T | \mathbf{X}) \gamma_B^{S_0}$ . Given  $\hat{\beta}_0, \hat{\gamma}_B^{S_0}, \hat{E}(\mathbf{B}|\mathbf{X})$  and  $\hat{\text{Var}}(\mathbf{B}|\mathbf{X})$ , we can easily obtain  $\hat{\gamma}_0^{S_k}$  by solving the equation  $E_{B|X}(\mu_0^{S_k}) - \text{logit}^{-1}(\hat{\beta}_0) = 0$ .

After obtaining  $\gamma_0^{S_k}$ , we then estimate  $\gamma_X^{S_k} = (\gamma_{X_1}^{S_k}, \dots, \gamma_{X_{P_k}}^{S_k})^T$  according to the following equation provided in Neuhaus and Jewell [1993]:

$$\beta_{X_p} = \left\{ \gamma_{X_p}^{S_k} + [E(\mathbf{B}^T | \mathbf{X} + 1_p) - E(\mathbf{B}^T | \mathbf{X})] \gamma_B^{S_0} \right\} \left\{ 1 - \frac{\text{Var}_{B|X}(\mu_0^{S_k})}{1 - E_{B|X}(\mu_0^{S_k})[1 - E_{B|X}(\mu_0^{S_k})]} \right\}$$

where  $1_p$  is a zero vector with the  $p^{\text{th}}$  term equals to 1 and  $p \in \{1, \dots, P_k\}$ . Similar to equation C.1, we can also obtain the Taylor-series-expansion estimation for  $E_{B|X}[(\mu_0^{S_k})^2] = E_{B|X}[\text{logit}^{-2}(\gamma_0^{S_k} + \mathbf{B}^T \gamma_B^{S_0})] \approx \frac{e^{2w}}{(1+e^w)^2} \left[ 1 + \frac{1}{2} \frac{2-e^w}{(1+e^w)^2} \sum_{i=1}^Q \sum_{j=1}^Q \gamma_{B_i}^{S_0} \gamma_{B_j}^{S_0} \text{Cov}(B_i, B_j | \mathbf{X}) \right]$ , together with  $E_{B|X}(\mu_0^{S_k})$ , we then obtain an approximation of  $V_{B|X}(\mu_0^{S_k}) = E_{B|X}[(\mu_0^{S_k})^2] - E_{B|X}(\mu_0^{S_k})^2$ . Given  $\hat{\beta}_X, \hat{\gamma}_B^{S_0}, \hat{E}(\mathbf{B}|\mathbf{X})$ , and  $\hat{\gamma}_0^{S_k}$ , we can obtain  $\hat{\gamma}_X^{S_k} = \hat{\beta}_X (1 - \frac{\hat{V}_{B|X}(\mu_0^{S_k})}{\hat{E}_{B|X}(\mu_0^{S_k})[1 - \hat{E}_{B|X}(\mu_0^{S_k})]})^{-1} - [\hat{E}(\mathbf{B}^T | \mathbf{X} + 1_p) - \hat{E}(\mathbf{B}^T | \mathbf{X})] \hat{\gamma}_B^{S_0}$ .

Note that we estimate  $E(\mathbf{B}|\mathbf{X}; \boldsymbol{\theta}) = g'^{-1}(\boldsymbol{\theta}\mathbf{X})$  and  $\text{Var}(\mathbf{B}|\mathbf{X}; \boldsymbol{\theta}) = g'^{-1}(\boldsymbol{\theta}\mathbf{X}) [1 - g'^{-1}(\boldsymbol{\theta}\mathbf{X})]$  using the internal data by regressing each  $\mathbf{B}$  on  $\mathbf{X}$  with appropriate link function  $g'()$  based on the type of  $\mathbf{B}$ , e.g., when  $\mathbf{B}$  is continuous, linear regression and identity link is used; when  $\mathbf{B}$  is binary, logistic regression and logit link is used. Given  $\hat{\boldsymbol{\theta}}, \hat{E}(\mathbf{B}|\mathbf{X}) = \hat{\boldsymbol{\theta}}^T E(\mathbf{X})$  is used.

## C.2 Additional Simulation Results

In this section, we show the results of additional simulations to assess the performance of the proposed strategy for point estimates and variance estimation.

### C.2.1 Continuous outcome Y (a supplement to Simulation I in Chapter 4)

**Goal:** To examine the proposed method when the outcome is continuous and the target model is linear regression.

**Simulation setup:** This simulation is the same as Simulation I in Chapter 4, except the generative outcome model now follows Gaussian distribution:

$$\begin{cases} \text{Internal:} & Y|\mathbf{X}, \mathbf{B} \sim N(-1 - X_1 - X_2 - B_1 - B_2, 1); \\ \text{External 1:} & Y|\mathbf{X}, \mathbf{B} \sim N(1 - X_1 - X_2 - B_1 - B_2, 1); \\ \text{External 2:} & Y|\mathbf{X}, \mathbf{B} \sim N(3 - X_1 - X_2 - B_1 - B_2, 1). \end{cases}$$

The target outcome model (model 2 in Chapter 4) is now a linear regression:

$$E(Y|\mathbf{X}, \mathbf{B}, \mathbf{S}) = \gamma_0 + \sum_{k=1}^2 \gamma_0^{S_k} S_k + \sum_{p=1}^2 \gamma_{X_p} X_p + \sum_{q=1}^2 \gamma_{B_q} B_q,$$

**Results:** Figure C.1 shows similar pattern as those in Simulation I in Chapter 4, where the proposed method (red dotted curve) has the smallest bias among all for all covariates (Figure C.1a), largest precision gain compared with others (Figure C.1b), and the closest variance estimation to the Monte Carlo empirical variance (Figure C.1c).

### C.2.2 Smaller covariate effect (a modification to Simulation I in Chapter 4)

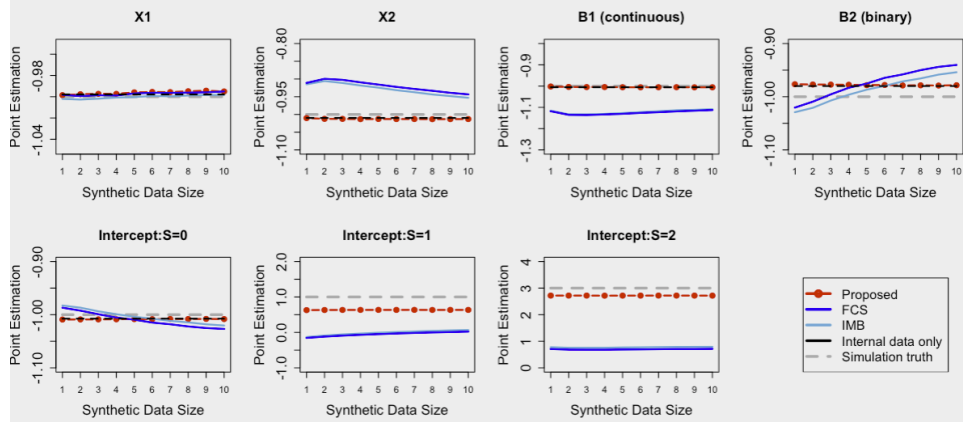
**Goal:** To assess our approach when the magnitude and the difference of covariate effects are small across different populations in the target outcome model.

**Simulation setup:** This simulation is the same as Simulation I in Chapter 4, except the coefficient effect is now -0.5 instead of -1, and the intercept difference is smaller among populations:

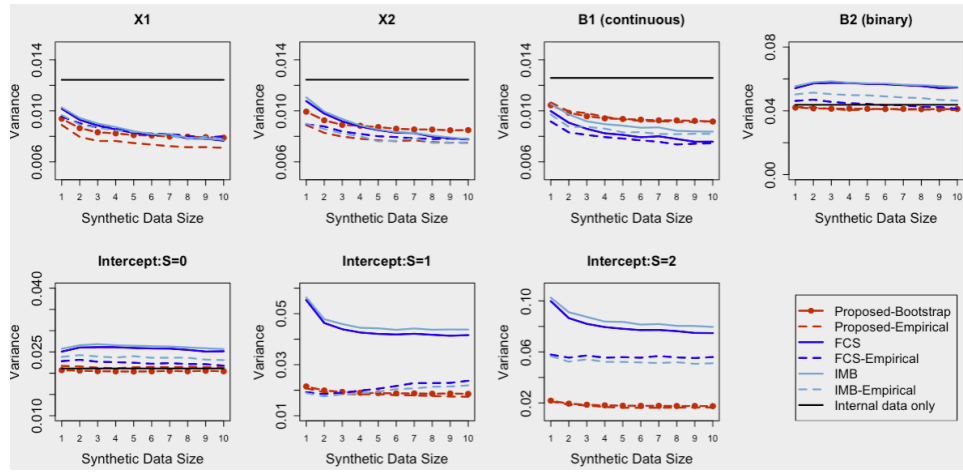
$$\begin{cases} \text{Internal:} & \text{logit}[\Pr(Y = 1|\mathbf{X}, \mathbf{B})] = -1 - 0.5(X_1 + X_2 + B_1 + B_2), \text{ prevalence}=0.28; \\ \text{External 1:} & \text{logit}[\Pr(Y = 1|\mathbf{X}, \mathbf{B})] = -0.5 - 0.5(X_1 + X_2 + B_1 + B_2), \text{ prevalence}=0.36; \\ \text{External 2:} & \text{logit}[\Pr(Y = 1|\mathbf{X}, \mathbf{B})] = 0 - 0.5(X_1 + X_2 + B_1 + B_2), \text{ prevalence}=0.45. \end{cases}$$

**Results:** Figure C.2a shows that compared with larger covariate effects in Simulation I in Chapter 4, when the X covariate effect is small, FCS and IMB have smaller bias in estimating X coefficients but still lack the ability to identify population-specific effects (i.e. intercepts of external populations). Similarly, Figure C.2b shows smaller bias of variance estimation. Note that the Rubin's rule variance estimator in Figure C.2b is too large (the pink curve) so that it falls outside of the range of the figure.

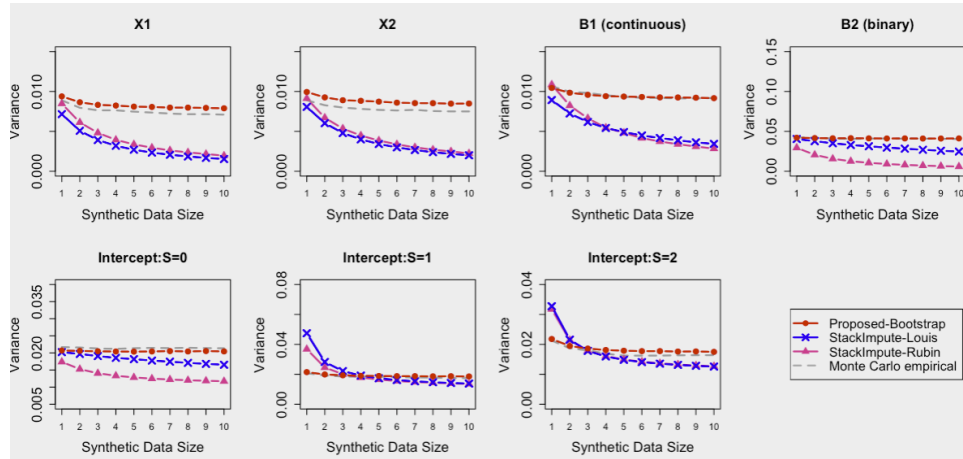




(a) Point estimates

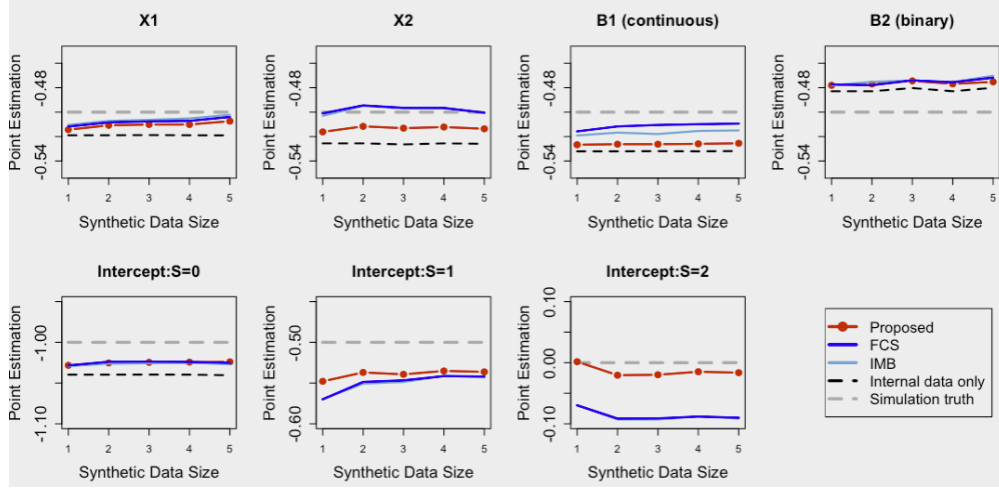


(b) Variance estimator vs. the empirical variance

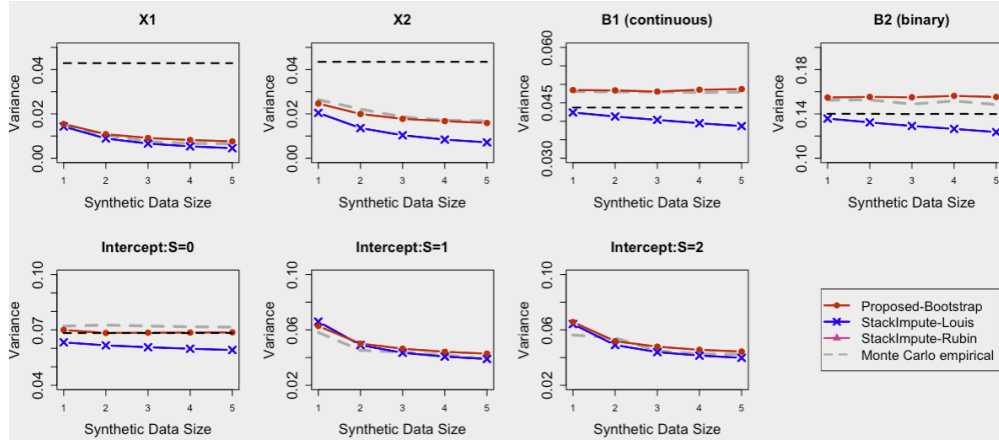


(c) Different variance estimators of the proposed method

Figure C.1: Results of Simulation C.2.1 over increasing synthetic data size (a) point estimates (b) variance estimation vs. Monte Carlo empirical variance (c) different variance estimators of the proposed method.



(a) Point estimates



(b) Different variance estimators of the proposed method

Figure C.2: Results of Simulation C.2.2 over increasing synthetic data size (a) point estimates (b) different variance estimators of the proposed method.

### C.2.3 Different X covariate effects in the outcome model (a more flexible outcome model compared with Simulation I in Chapter 4)

**Goal:** In Chapter 4, we only present the simulation results allowing the target model's intercept to differ across populations. In this simulation, we additionally show the performance of the proposed method when all possible X covariates coefficients are allowed to differ across populations (similar to model 1 or “different intercept and covariates” model in the real data example in Chapter 4).

**Simulation setup:** This simulation is the same as Simulation I in Chapter 4 except now that the

generative outcome models are as follows:

$$\begin{cases} \text{Internal:} & \text{logit}[\Pr(Y = 1|\mathbf{X}, \mathbf{B})] = -1 - X_1 - X_2 - B_1 - B_2, \text{ prevalence}= 0.3; \\ \text{External 1:} & \text{logit}[\Pr(Y = 1|\mathbf{X}, \mathbf{B})] = 1 + X_1 - X_2 - B_1 - B_2, \text{ prevalence}= 0.58; \\ \text{External 2:} & \text{logit}[\Pr(Y = 1|\mathbf{X}, \mathbf{B})] = 3 + 3X_1 + 3X_2 - B_1 - B_2, \text{ prevalence}= 0.70. \end{cases}$$

**Results:** Similar to the results of Simulation I in Chapter 4, the results in Figure C.3 shows outstanding performance of the proposed method in both point estimates and variance estimation compared with others. For example, the proposed method has small bias less than 0.02 when estimating  $X_2$  in population  $S=1$  while the bias in FCS and IMB can go up to 0.78 (i.e. almost 40 times of the proposed method).

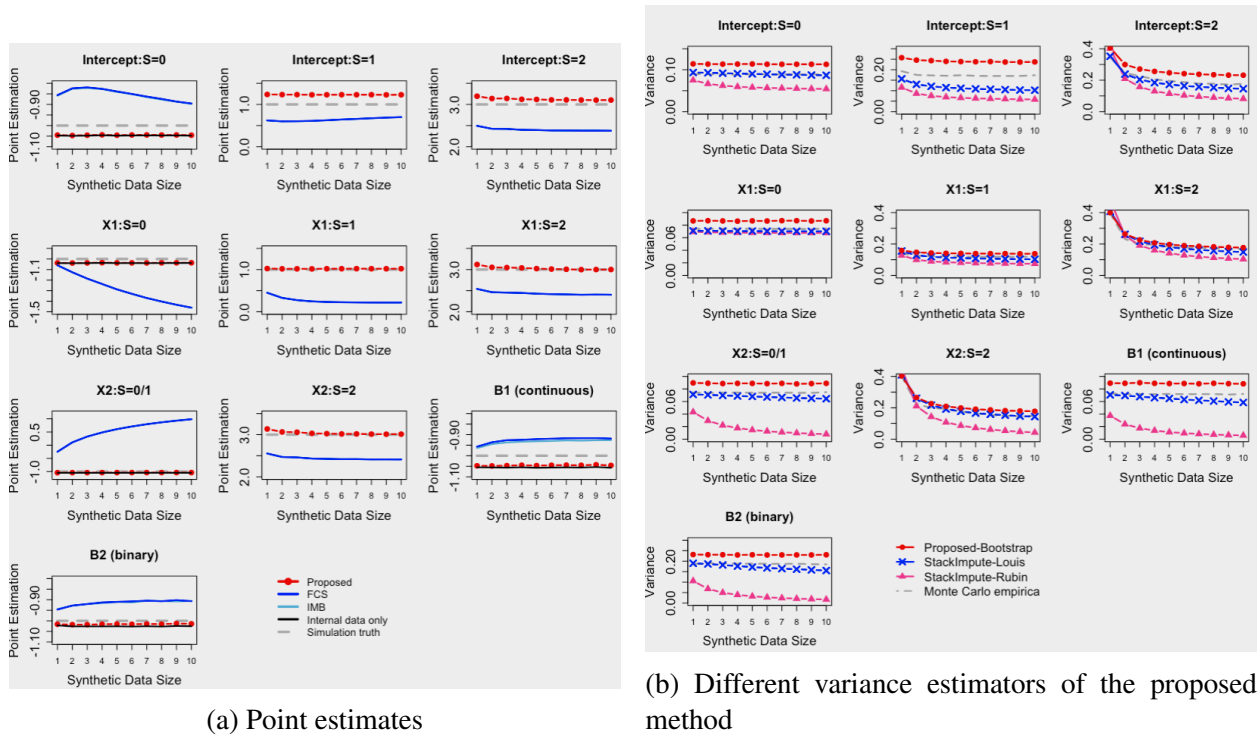


Figure C.3: Results of Simulation C.2.3 over increasing synthetic data size (a) point estimates (b) different variance estimators of the proposed method.

## C.2.4 Violation of transportability assumption

**Goal:** To examine the proposed method when Assumption 2 ( $X_{\text{miss}}|X_{\text{obs}}$  and  $B|X$  are transportable between the internal and the external populations) is violated. We present two examples where the violation only causes ignorable bias in case 1 while it has larger impact in case 2.

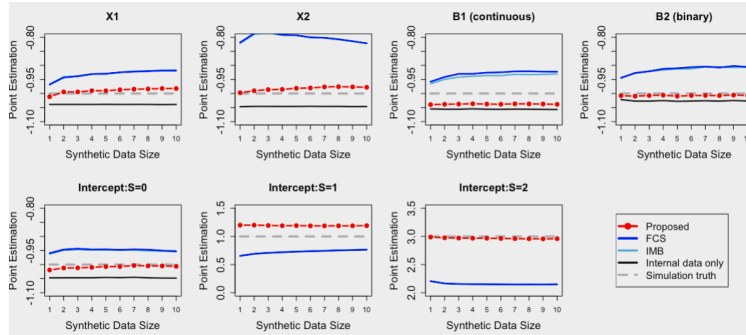
### C.2.4.1 Case 1: different $B|X$ distribution in external population 2

**Simulation setup:** This simulation is the same as Simulation I in Chapter 4 except that now the external model 2 has different marginal  $B_1$  distribution and different conditional distribution  $B_2|X_1, X_2, B_1$ :

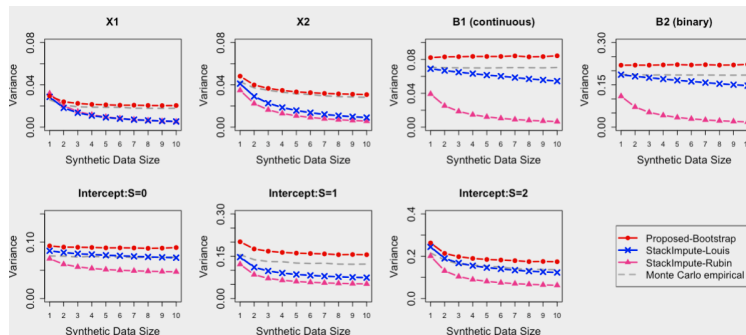
- $B_1$  has mean 1.5 and standard deviation 1.5 in external population 2 while in other populations  $B_1$  has mean 0 and standard deviation 1;
- $B_2|X, B_1 \sim \text{Ber}\{[1 + \exp^{-1}(0.2X_1 + 0.3X_2 + 0.4B_1)]\}$  in external population 2 while in other populations  $B_2|X, B_1 \sim \text{Ber}\{[1 + \exp^{-1}(0.1X_1 + 0.2X_2 + 0.3B_1)]\}$ .

Note that both  $B_1$  and  $B_2$  are only observed in the internal study and multiple imputations are needed for them, where  $B_2|X$  and  $B_2|X, B_1$  should be the same across populations according to Assumption 2.

**Results:** Figure C.4a indicates that the violation of transportability assumption in the proposed method has limited impact of point estimation with ignorable bias while Figure C.4b shows similar pattern of variance estimations as before.



(a) Point estimates



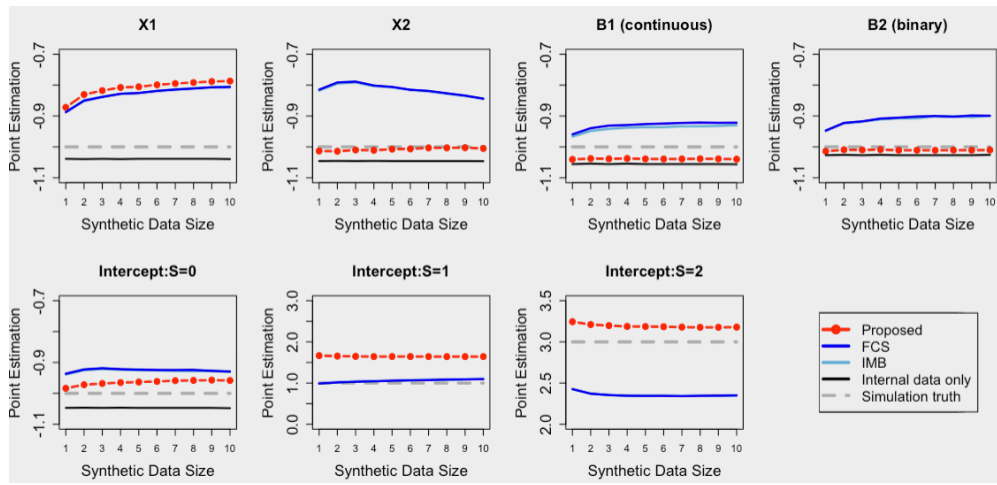
(b) Different variance estimators of the proposed method

Figure C.4: Results of Simulation C.2.4.1 over increasing synthetic data size (a) point estimates (b) different variance estimators of the proposed method.

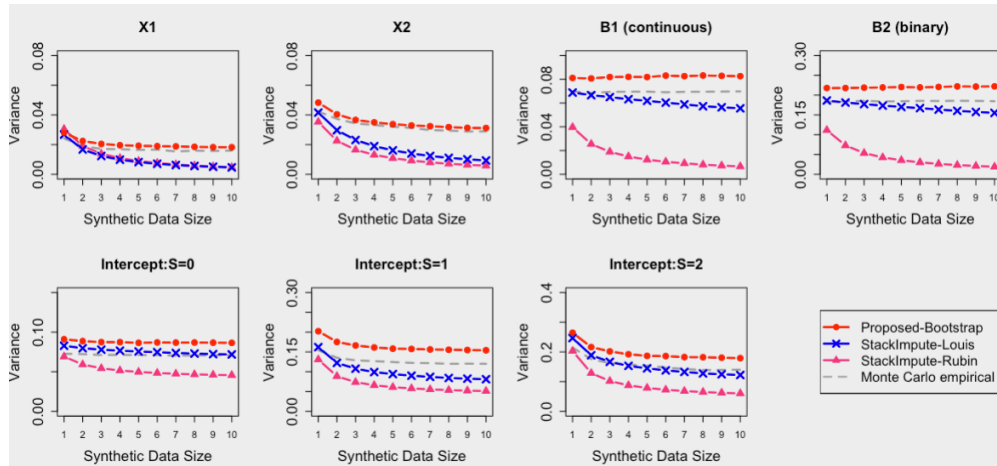
### C.2.4.2 Case 2: different marginal $X_1$ distribution in external populations

**Simulation setup:** This simulation is the same as Simulation I in Chapter 4 except now that in the external studies,  $X_1 \sim N(1, 1.5)$  while in the internal study  $X_1 \sim N(0, 1)$ . This will lead to different distribution conditional on  $X_1$  and thus violates Assumption 2.

**Results:** Figure C.5a shows that such violation leads to some bias of estimated coefficient  $X_1$ , i.e., 0.2 absolute bias. Besides that, the proposed method has nearly unbiased point estimates for other parameter (i.e. up to 0.014 absolute bias) while the bias in FCS and IMB can be up to 15 times the bias of the proposed method. Similarly, Figure C.5b shown unbiased variance estimation of the proposed bootstrap estimator.



(a) Point estimates



(b) Different variance estimators of the proposed method

Figure C.5: Results of Simulation C.2.4.2 over increasing synthetic data size (a) point estimates (b) different variance estimators of the proposed method.

## BIBLIOGRAPHY

- K. Viele, S. Berry, B. Neuenschwander, B. Amzal, F. Chen, N. Enas, B. Hobbs, J. G. Ibrahim, N. Kinnersley, S. Lindborg, S. Micallef, S. Roychoudhury, and L. Thompson. Use of historical control data for assessing treatment effects in clinical trials. *Pharmaceutical Statistics*, 13:41–54, 2014.
- D. Dejardin, P. Delmar, C. Warne, K. Patel, J. van Rosmalen, and E. Lesaffre. Use of a historical control group in a noninferiority trial assessing a new antibacterial treatment: A case study and discussion of practical implementation aspects. *Pharmaceutical Statistics*, 17:169–181, 2018.
- X. H. Li and Y. Song. Target population statistical inference with data integration across multiple sources—an approach to mitigate information shortage in rare disease clinical trials. *Statistics in Biopharmaceutical Research*, 12:322–333, 2020.
- C. Bycroft. Integrated household surveys: a survey vehicles approach. Technical report, Wellington, 1 2011.
- S. Yang and J. K. Kim. Statistical data integration in survey sampling: a review. *Japanese Journal of Statistics and Data Science*, 3:625–650, 2020.
- S. Yang and P. Ding. Combining multiple observational data sources to estimate causal effects. *Journal of the American Statistical Association*, 3:1540–1554, 2020a.
- M. H. Gail, L. A. Brinton, D. P. Byar, D. K. Corle, S. B. Green, C. Schairer, and J. J. Mulvihill. Projecting individualized probabilities of developing breast cancer for white females who are being examined annually. *Journal of the National Cancer Institute*, 81:1879–1886, 1989.
- I. M. Thompson, D. P. Ankerst, C. Chi, P. J. Goodman, C. M. Tangen, M. S. Lucia, Z. Feng, H. L. Parnes, and C. A. Coltman. Assessing prostate cancer risk: Results from the prostate cancer prevention trial. *European Urology*, 98:529–534, 2006.
- M. J. Roobol, H. A. van Vugt, S. Loeb, X. Zhu, M. Bul, C. H. Bangma, A. G. van Leenders, E. W. Steyerberg, and F. H. Schröder. Prediction of prostate cancer risk: the role of prostate volume and digital rectal examination in the ERSPC risk calculators. *European Urology*, 61:577–583, 2012.
- C. Stephan, B. Vogel, H. Cammann, M. Lein, V. Klevecka, P. Sinha, G. Kristiansen, D. Schnorr, K. Jung, and S. A. Loening. An artificial neural network as a tool in risk evaluation of prostate cancer. indication for biopsy with the psa range of 2-20 microg/l. *Der Urologe, Ausgabe A*, 42: 1221–1229, 2003. doi: <https://doi.org/10.1007/s00120-003-0322-7>.

- A. Osareh and B. Shadgar. Machine learning techniques to diagnose breast cancer. *2010 5th International Symposium on Health Informatics and Bioinformatics*, pp:114–120, 2010. doi:10.1109/HIBIT.2010.5478895.
- H. Estiri, Z.H. Strasser, J.G. Klann, P. Naseri, K. B. Waghlikar, and S. N. Murphy. Predicting covid-19 mortality with electronic medical records. *npj Digit. Med*, 4:15, 2021. doi: <https://doi-org.proxy.lib.umich.edu/10.1038/s41746-021-00383-x>.
- G. W. Imbens and T. Lancaster. Combining micro and macro data in microeconomic models. *The Review of Economic Studies*, 61:655–680, 1994.
- S. Grill, M. Fallah, R. J. Leach, I. M. Thompson, K. Hemminki, and D. P. Ankerst. A simple-to-use method incorporating genomic markers into prostate cancer risk prediction tools facilitated future validation. *Journal of Clinical Epidemiology*, 68:563–573, 2015.
- E. Bareinboim and J. Pearl. A general algorithm for deciding transportability of experimental results. *Journal of Causal Inference*, 1:107–134, 2013.
- J. Qin. Combining parametric and empirical likelihoods. *Biometrika*, 87:484–490, 2000.
- P. Han and J. F. Lawless. Empirical likelihood estimation using auxiliary summary information with different covariate distribution. *Statistics Sinica*, 29:1321–1342, 2019.
- N. Chatterjee, Y.-H. Chen, P. Maas, and R. J. Carroll. Constrained maximum likelihood estimation for model calibration using summary-level information from external big data sources. *Journal of the American Statistical Association*, 111:107–117, 2016.
- W. Cheng, J. M. G. Taylor, P. S. Vokonas, S. K. Park, and B. Mukherjee. Improving estimation and prediction in linear regression incorporating external information from an established reduced model. *Statistics in Medicine*, 37:1515–1530, 2018.
- W. Cheng, J. M. G. Taylor, T. Tomlins, and B. Mukherjee. Informing a risk prediction model for binary outcomes with external coefficient information. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 68:121–139, 2019.
- S. Grill, D. P. Ankerst, M. H. Gail, N. Chatterjee, and R. M. Pfeiffer. Comparison of approaches for incorporating new information into existing risk prediction models. *Statistics in Medicine*, 36:1134–1156, 2017.
- J. P. Estes, B. Mukherjee, and J. M. G. Taylor. Empirical Bayes estimation and prediction using summary-level information from external big data sources adjusting for violations of transportability. *Statistics in Biosciences*, 10:568–586, 2017.
- P. Kundu, R. Tang, and N. Chatterjee. Generalized meta-analysis for multiple regression models across studies with disparate covariate information. *Biometrika*, 106:567–585, 2019.
- Z. Chen, J. Ning, Y. Shen, and J. Qin. Combining primary cohort data with external aggregate information without assuming comparability. *Biometrics*, 2020. doi: 10.1111/biom.13356.

- S. Yang and P. Ding. Combining multiple observational data sources to estimate causal effects. *Journal of the American Statistical Association*, 115:1540–1554, 2020b.
- T. E. Raghunathan, J. P. Reiter, and D. B. Rubin. Multiple imputation for statistical disclosure limitation. *Journal of Official Statistics*, 19:1, 2003.
- J. P. Reiter and S. K. Kinney. Inferentially valid, partially synthetic data: Generating from posterior predictive distributions not necessary. *Journal of Official Statistics*, 28:583–n/a, 2012.
- J. P. Reiter. Satisfying disclosure restrictions with synthetic data sets. *Journal of Official Statistics*, 18:531, 2002.
- D. B. Rubin. *Multiple Imputation for Nonresponse in Surveys*. John Wiley & Sons Inc., New York, 1987.
- S. Van Buuren and C. G. M. Oudshoorn. MICE: Multivariate imputation by chained equations in R. *Journal of Statistical Software*, 45, 2011. doi: 10.18637/jss.v045.i03.
- D. Xu, M. J. Daniels, and A. G. Winterstein. Sequential BART for imputation of missing covariates. *Biostatistics*, 17:589–602, 2016.
- S. A. Tomlins, J. R. Day, R. J. Lonigro, D. H. Hovelson, J. Siddiqui, L. P. Kunju, R. L. Dunn, S. Meyer, P. Hodge, J. Groskopf, J. T. Wei, and A. M. Chinnaiyan. Urine TMPRSS2:ERG plus PCA3 for individualized prostate cancer risk assessment. *European Urology*, 70:45–53, 2015.
- H. Zhang, L. Deng, M. Schiffman, J. Qin, and K. Yu. Generalized integration model for improved statistical inference by leveraging external summary data. *Biometrika*, 107:689–703, 2020.
- P. Goos and B. Jones. *Optimal Design of Experiments: A Case Study Approach*. Wiley, New Jersey, 2011.
- M. Truong, B. Yang, and D. F. Jarrard. Toward the detection of prostate cancer in urine: a critical analysis. *Journal of Urology*, 189:422–429, 2013.
- A. M. Bohnen, F. P. Groeneveld, and J. L. H. R. Bosch. Serum prostate-specific antigen as a predictor of prostate volume in the community: the krimpen study. *European Urology*, 51:1645–1653, 2007.
- A. J. Simpkin, J. L. Donovan, K. Tilling, J. Athene Lane, R. M. Martin, P.C. Albertsen, A. Bill-Axelson, H. Ballentine Carter, J. L. Bosch, L. Ferrucci, F. C. Hamdy, L. Holmberg, E. Jeffrey Metter, D. E. Neal, C. C. Parker, and C. Metcalfe. Prostate-specific antigen patterns in us and european populations: comparison of six diverse cohorts. *BJU International*, 118:911–918, 2016.
- R. E. Schapire. The strength of weak learnability. *Machine Learning*, 5:197–227, 1990.
- Y. Freund and R. R. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55:119–139, 1997.
- L. Breiman. Bagging predictors. *Machine Learning*, 24:123–140, 1996.



- L. Breiman. Bagging predictors. *Random forests*, 45:5–32, 2001.
- M. J. van der Laan, E. C. Polley, and A. E. Hubbard. Super learner. *Statistical Applications in Genetics and Molecular Biology*, 6:25, 2007.
- F. Wang, L. Wang, and P. X.-K. Song. Quadratic inference function approach to merging longitudinal studies: validation and joint estimation. *Biometrika*, 99(3):755–762, 09 2012a. ISSN 0006-3444. doi: 10.1093/biomet/ass021.
- J. Antonelli, C. Zigler, and F. Dominici. Guided Bayesian imputation to adjust for confounding when combining heterogeneous data sources in comparative effectiveness research. *Biostatistics*, 18:553–568, 2017.
- C. Wang, G. Parmigiani, and F. Dominici. Bayesian effect estimation accounting for adjustment uncertainty. *Biometrics*, 68:661–671, 2012b.
- L. J. Beesley and J. M. G. Taylor. A stacked approach for chained equations multiple imputation incorporating the substantive model. *Biometrics*, 2020. doi: 10.1111/biom.13372.
- L. J. Beesley and J. M. G. Taylor. Accounting for not-at-random missingness through imputation stacking. *ArXiv*, 2021. doi: 2101.07954.
- J. Neuhaus and N. Jewell. A geometric approach to assess bias due to omitted covariates in generalized linear models. *Biometrika*, 80:807–815, 1993.
- S. Van Buuren, J. P. L. Brand, C. G. M. Groothuis-Oudshoorn, and D. B. Rubin. Fully conditional specification in multivariate imputation. *Journal of Statistical Computation and Simulation*, 76: 1049–1064, 2006.
- F. Li, M. Baccini, F. Mealli, E. R. Zell, C. E. Frangakis, and D. B. Rubin. Multiple imputation by ordered monotone blocks with application to the anthrax vaccine research program. *Journal of Computational and Graphical Statistics*, 23(3):877–892, 2014. doi: 10.1080/10618600.2013.826583.
- R. J. A. Little and D. B. Rubin. *Statistical Analysis with Missing Data*. John Wiley and Sons, Inc, Hoboken, NJ, 2nd edition, 2002.
- J. Y. Dai, C. Kooperberg, M. Leblanc, and R. L. Prentice. Two-stage testing procedures with independent filtering for genome-wide gene-environment interaction. *Biometrika*, 99:929–944, 2012.
- C. Huang and J. Qin. A unified approach for synthesizing population-level covariate effect information in semiparametric estimation with survival data. *Statistics in Medicine*, 39:1573–1590, 2020.
- S. Rasser. Data fusion: identification problems, validity, and multiple imputation. *Statistica Sinica*, 33:153–171, 2004.
- J. P. Reiter. Bayesian finite population imputation for data fusion. *Statistica Sinica*, 22:795–811, 2012.

- C. Gouriéroux and A. Monfort. On the problem of missing data in linear models. *Review of Economic Studies*, 48:579–586, 1981.
- G. H. Golub and C. F. Van Loan. *Matrix Computations*. Johns Hopkins University, Baltimore, MD, 3rd edition, 1996.